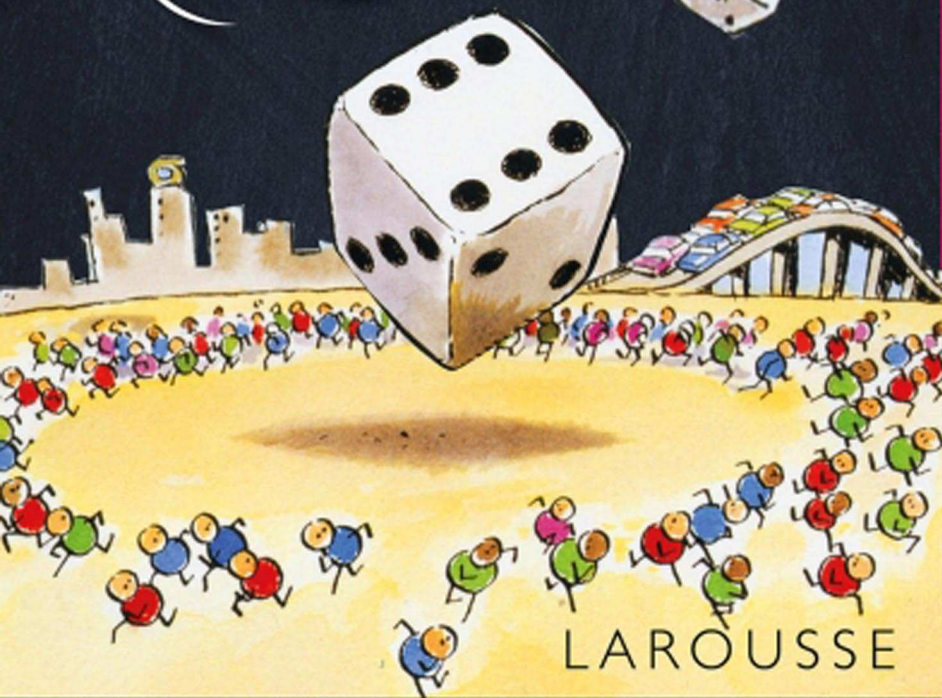


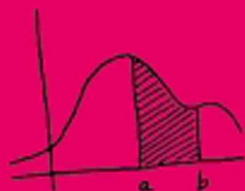
LARRY GONICK
ART HUFFMAN

LES STATISTIQUES EN BD

GRÂCE
À LARRY GONICK,
LES SCIENCES RETROUVENT
LE SOURIRE ET LES CONCEPTS
DIFFICILES DEVIENNENT
LIMPIDES.



LAROUSSE



$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad \text{ou} \quad \sum_{i=1}^n \frac{x_i}{n}$$

LES STATISTIQUES EN BD

AUTRES OUVRAGES DE LARRY GONICK :

AUX ÉDITIONS LAROUSSE (EN FRANÇAIS) :

LA CHIMIE EN BD
LES MATHS EN BD
LA GÉNÉTIQUE EN BD
LA PHYSIQUE EN BD

ET EN ANGLAIS :

THE CARTOON HISTORY OF THE UNIVERSE, VOLUMES 1-7
THE CARTOON HISTORY OF THE UNIVERSE II, VOLUMES 8-13
THE CARTOON HISTORY OF THE UNIVERSE III, VOLUMES 14-19
THE CARTOON HISTORY OF THE UNITED STATES
THE CARTOON GUIDE TO THE COMPUTER
THE CARTOON GUIDE TO THE ENVIRONMENT (AVEC ALICE OUTWATER)
THE CARTOON GUIDE TO (NON) COMMUNICATION
THE CARTOON GUIDE TO SEX (AVEC CHRISTINE DEVALT)

LES STATISTIQUES EN BD



LARRY GONICK
& WOOLLCOTT SMITH

Édition originale :

Ce titre a été publié pour la première fois en 1993 sous le titre original
The Cartoon Guide to Statistics © 1993 by Larry Gonick and Woolcott Smith.
Published by arrangement with William Morrow, an imprint of HARPERCOLLINS PUBLISHERS.
All rights reserved.

Édition française :

Éditions Larousse © 2016 pour l'édition en langue française

Direction de la publication : Carine Girac-Marinier

Direction éditoriale : Claude Nimmo

Direction éditoriale adjointe : Julie Pelpel-Moulian

Édition : Antoine Caron

Réalisation et coordination éditoriale : Belle Page

Traduction : Thierry Lafay, maître de conférences à l'université Paris 1 Panthéon-Sorbonne. Il enseigne notamment les mathématiques, les statistiques et les techniques quantitatives à l'École de Management de la Sorbonne.

© Larousse 2016

Toute reproduction ou représentation intégrale ou partielle, par quelque procédé que ce soit, de la nomenclature et/ou du texte et des illustrations contenus dans le présent ouvrage, et qui sont la propriété de l'Éditeur, est strictement interdite. Les Éditions Larousse utilisent des papiers composés de fibres naturelles, renouvelables, recyclables et fabriquées à partir de bois issus de forêts qui adoptent un système d'aménagement durable. En outre, les Éditions Larousse attendent de leurs fournisseurs de papier qu'ils s'inscrivent dans une démarche de certification environnementale reconnue.

ISBN : 978-2-03-593291-4

Sommaire

CHAPITRE 1.....	1
QU'EST-CE QUE LA STATISTIQUE?	
CHAPITRE 2.....	7
STATISTIQUES DESCRIPTIVES	
CHAPITRE 3.....	27
LES PROBABILITÉS	
CHAPITRE 4.....	53
LES VARIABLES ALÉATOIRES	
CHAPITRE 5.....	73
UNE HISTOIRE DE DEUX DISTRIBUTIONS	
CHAPITRE 6.....	89
ÉCHANTILLONNAGE	
CHAPITRE 7.....	111
INTERVALLES DE CONFIANCE	
CHAPITRE 8.....	137
TESTS D'HYPOTHÈSES	
CHAPITRE 9.....	157
COMPARAISON DE DEUX POPULATIONS	
CHAPITRE 10.....	181
MÉTHODES EXPÉRIMENTALES	
CHAPITRE 11.....	187
RÉGRESSION LINÉAIRE	
CHAPITRE 12.....	211
CONCLUSION	
BIBLIOGRAPHIE.....	221
INDEX.....	224

À PROPOS DES AUTEURS

WOOLCOTT SMITH EST PROFESSEUR ÉMÉRITE DE STATISTIQUES À L'UNIVERSITÉ DE TEMPLE. TITULAIRE D'UNE LICENCE ET D'UN MASTÈRE DE L'UNIVERSITÉ DU MICHIGAN ET D'UNE THÈSE DE L'UNIVERSITÉ JOHNS- HOPKINS, IL EST AUTEUR ET COAUTEUR DE PLUS D'UNE QUARANTAINE DE PUBLICATIONS DANS DIFFÉRENTS DOMAINES TELS QUE LES MARÉES NOIRES, LA THÉORIE STATISTIQUE ET LES STATISTIQUES ENVIRONNEMENTALES. IL A CONSEILLÉ DE NOMBREUX PROGRAMMES SCIENTIFIQUES NATIONAUX. IL ENTRETIENT DES DÉBATS ANIMÉS ET FAIT DU KAYAK AVEC SA FEMME LEAH ET SES DEUX ENFANTS MAINTENANT ADULTES KESTON ET AMELIA.



LARRY GONICK EST L'AUTEUR ET LE COAUTEUR DE GUIDES ILLUSTRÉS ET DE BANDES DESSINÉES AYANT REÇU LE PRIX HARVEY. IL A ÉCRIT ET DESSINÉ LES CLASSIQUES DE LA SCIENCE DANS *DISCOVER MAGAZINE* ET A CRÉÉ LA BANDE DESSINÉE *KOKOPELLI & COMPANY* POUR LE MAGAZINE *MUSE*. IL EST L'UNIQUE DESSINATEUR À AVOIR CONTRIBUÉ AVEC DES DESSINS COMIQUES ORIGINAUX À LA PRESTIGIEUSE REVUE *SCIENCE*. APRÈS AVOIR ABANDONNÉ LES MATHÉMATIQUES À LA HARVARD SCHOOL, IL S'EST INSTALLÉ SUR LA CÔTE PACIFIQUE AVEC SA FEMME LISA. LEURS DEUX FILLES SOPHIE ET ANNA ONT DÉSORMAIS QUITTÉ LE COCON FAMILIAL.

Remerciements

NOUS SERONS ÉTERNELLEMENT RECONNAISSANTS ENVERS NOTRE ÉDITEUR ORIGINAL CAROL COHEN DE HAPPERCOLLINS QUI NOUS A SUGGÉRÉ CE PROJET ET ENVERS NOTRE AGENTE LITTÉRAIRE VICKY BIJUR QUI NOUS A PERMIS DE LE MENER À BIEN GRÂCE AUX COAUTEURS. UN COUP DE CHAPEAU PARTICULIER AUX ÉDITEURS ET ASSISTANTS AVEC LESQUELS NOUS AVONS TRAVAILLÉ DEPUIS : ERICA SPABERG, STEPHANIE MEYER, PETER HUBBARD, COLE HAGER ET NICK AMPHLETT.

LES COMMENTAIRES DE WILLIAM FAIRLEY ET LEAH SMITH ONT CONTRIBUÉ À L'AMÉLIORATION DES VERSIONS PRÉLIMINAIRES.

DONNA OKINO A FOURNI UNE ASSISTANCE ET DES CONSEILS INESTIMABLES POUR LA PRODUCTION ET LA PRÉPARATION DES PAGES DESSINÉES. ELLE DIT QUE CRÉER UN GUIDE ILLUSTRÉ EST PLUS DIFFICILE QUE DE COURIR UN MARATHON ET ELLE SAIT DE QUOI ELLE PARLE : ELLE A FAIT LES DEUX.

DES REMERCIEMENTS PARTICULIERS AUX CRÉATEURS DE FONTOGRAPHER, LE LOGICIEL MAGIQUE QUI NOUS A PERMIS DE CRÉER PLUSIEURS POLICES DE CARACTÈRES QUI SIMULENT L'ÉCRITURE MANUELLE.

ET COMME L'ÉDUCATION N'EST PAS À SENS UNIQUE, UN COUP DE CHAPEAU À LA PATIENCE DES ÉTUDIANTS DE L'UNIVERSITÉ DE TEMPLE. LE FUTUR LEUR APPARTIENT.

NOTE DU TRADUCTEUR : DANS CE LIVRE, NOUS AVONS CONSERVÉ LES DONNÉES DU MANUSCRIT ORIGINAL BIEN QUE LES UNITÉS SOIENT AMÉRICAINES. EN EFFET, LES STATISTIQUES TRAVAILLENT SUR LES CHIFFRES ET PEU IMPORTE L'UNITÉ CHOISIE. IL NE NOUS SEMBLAIT DONC PAS NÉCESSAIRE DE CONVERTIR TOUS CES JEUX DE DONNÉES, CHOISIS PAR LES AUTEURS POUR LEUR INTÉRÊT STATISTIQUE ET ILLUSTRÉS PAR DES DESSINS ORIGINAUX. NOUS AVONS TOUTEFOIS DONNÉ LE RATIO POUR CONVERTIR CES CHIFFRES EN UNITÉS FRANÇAISES, LES CONCLUSIONS ÉTANT AUSSI AGRÉMENTÉES DE LEUR CONVERSION.



Chapitre I

Qu'est-ce que la statistique?

ON SE DÉBROUILLE DANS LA VIE EN FAISANT DES CHOIX FONDÉS
SUR UNE INFORMATION INCOMPLÈTE...



LA PLUPART D'ENTRE NOUS VIVENT
SANS PROBLÈME AVEC CE NIVEAU
D'INCERTITUDE.



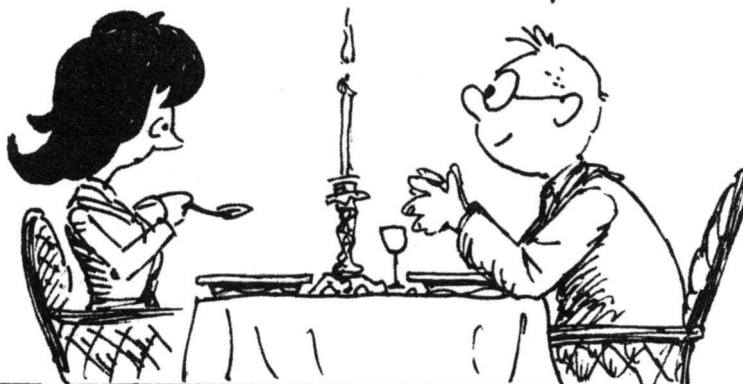
JE VAIS PRENDRE
DE LA SOUPE
S'IL VOUS
PLÂT.

AAAAH... POURRIEZ-VOUS
M'APPORTER D'ABORD
UNE **CALCULATRICE** ?



CE QUI REND LES STATISTIQUES SI UNIQUES, C'EST LA POSSIBILITÉ
DE **QUANTIFIER** L'INCERTITUDE, DE FAÇON À LA RENDRE PLUS PRÉCISE.
LES STATISTICIENS PEUVENT ALORS FAIRE AVEC ASSURANCE
DES **PROPOSITIONS CATÉGORIQUES** SUR LEUR NIVEAU D'INCERTITUDE !

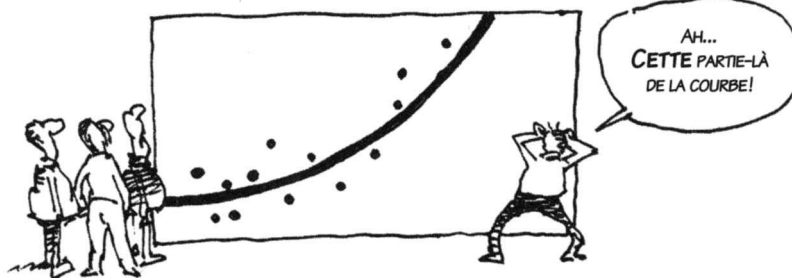
TRÈS BON CHOIX ! JE SUIS SÛR
À 95 % QUE LA SOUPE DE CE **SOIR**
A UNE PROBABILITÉ ENTRE 73 ET 77 %
D'ÊTRE **VRAIMENT DÉLICIEUSE** !



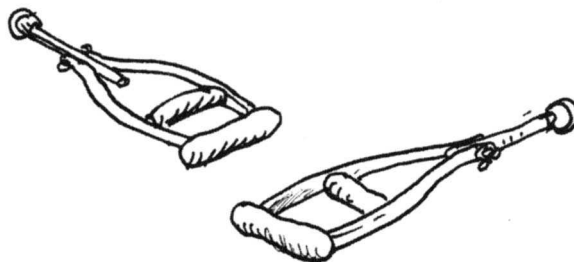
CELA NE SERT PAS UNIQUEMENT
À COMMANDER UNE SOUPE!
LA STATISTIQUE CONCERNE
AUSSI DES QUESTIONS
DE VIE OU DE MORT...



PAR EXEMPLE, EN 1986, LA NAVETTE SPATIALE **CHALLENGER** A EXPLODÉ, TUANT SEPT ASTRONAUTES. LA DÉCISION DE LANCER LA NAVETTE ALORS QU'IL FAISAIT -1°C A ÉTÉ PRISE SANS QU'AUCUNE ANALYSE SIMPLE AIT ÉTÉ FAITE SUR LES PERFORMANCES DE LA NAVETTE À BASSE TEMPÉRATURE.



LE **VACCIN DE SALK** CONTRE LA POLIOMYÉLITE FOURNIT UN EXEMPLE PLUS POSITIF. EN 1954, DES ESSAIS CLINIQUES DU VACCIN FURENT TESTÉS SUR PLUS DE 400 000 ENFANTS, AVEC UN CONTRÔLE STRICT POUR ÉLIMINER DES RÉSULTATS BIAISÉS. LES ANALYSES STATISTIQUES DES RÉSULTATS ONT ÉTÉ CONCLUANTES. ELLES ONT PERMIS D'ÉTABLIR L'EFFICACITÉ DU VACCIN ET, AUJOURD'HUI, LA POLIOMYÉLITE EST PRESQUE ÉRADIQUÉE.



POUR ACCOMPLIR LEURS EXPLOITS DE DÉTECTIVE DU MONDE RÉEL,
LES STATISTICIENS UTILISENT TROIS DISCIPLINES LIÉES :

L'analyse des données

LA COLLECTE, LA PRÉSENTATION
ET LE RÉSUMÉ DES DONNÉES.

Les probabilités

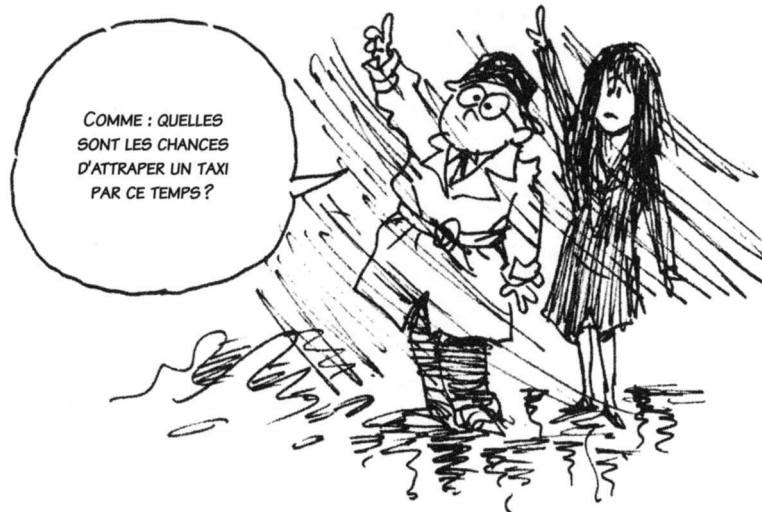
LES LOIS DU HASARD, QUE CE SOIT
DANS OU EN DEHORS D'UN CASINO.

L'inférence statistique

LA SCIENCE QUI CONSISTE À ÉLABORER
DES CONCLUSIONS STATISTIQUES À PARTIR
DE DONNÉES SPÉCIFIQUES EN UTILISANT
LES CONNAISSANCES SUR LES PROBABILITÉS.



DANS CET OUVRAGE, NOUS ALLONS ABORDER CES TROIS DISCIPLINES EN LES APPLIQUANT
À DE NOMBREUX TYPES DE SITUATIONS OÙ LES STATISTIQUES JOUENT UN RÔLE CRUCIAL
DANS LE MONDE MODERNE.



DANS LE CHAPITRE 2,
NOUS EXAMINERONS UN ENSEMBLE
SIMPLE DE DONNÉES : LES POIDS
D'UN GROUPE D'ÉTUDIANTS DE LICENCE.



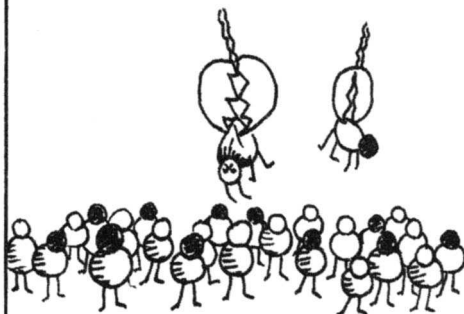
DANS LE CHAPITRE 3, NOUS ÉTUDIERONS
LES LOIS DE PROBABILITÉS LÀ OÙ ELLES
SONT APPARUES, LES MAISONS DE JEUX.



LES CHAPITRES 4 ET 5 MONTRERONT
COMMENT DÉCRIRE LE MONDE GRÂCE
AUX MODÈLES DE PROBABILITÉS,
ET AU CONCEPT DE VARIABLE ALÉATOIRE.



LE CHAPITRE 6 INTRODUIRA L'UNE
DES PREMIÈRES PROCÉDURES DU BON
STATISTICIEN : CONSTRUIRE UN ÉCHANTILLON
À PARTIR D'UNE GRANDE POPULATION.

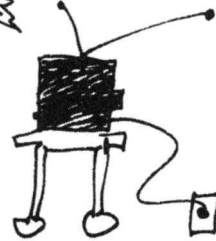


DANS LE CHAPITRE 7
ET LES SUIVANTS, NOUS DÉCRIRONS
COMMENT FAIRE DES INFÉRENCES
STATISTIQUES DANS DES SITUATIONS
RÉELLES AUSSI COURANTES
QUE LES PRÉVISIONS ÉLECTORALES,
LE CONTRÔLE DE QUALITÉ
INDUSTRIELLE, LES TESTS MÉDICAUX,
LE SUIVI ENVIRONNEMENTAL,
LES DISCRIMINATIONS ET LE DROIT.



ENFIN, LORSQUE NOUS PARLONS DE STATISTIQUES IL EST DIFFICILE D'ÉVITER DE MENTIONNER LA **MÉFIANCE** TRÈS RÉPANDUE À L'ÉGARD DES STATISTIQUES AUJOURD'HUI. ON A TOUS ENTENDU DIRE QUE LES STATISTIQUES SONT LA « FORME LA PLUS ÉLABORÉE DU MENSONGE », ET IL EST PRESQUE IMPOSSIBLE DE TROUVER DE BONNES ANALYSES STATISTIQUES DANS LA VIE DE TOUS LES JOURS. ALORS QUE FAIRE ?

3 DOCTEURS SUR 4 RECOMMANDENT DE NE PAS CROIRE UNE PROPOSITION QUI COMMENCE PAR « 3 DOCTEURS SUR 4 »...



NOTRE HUMBLE OPINION EST QU'EN **APPRENDRE UN PEU PLUS SUR LE SUJET** N'EST SÛREMENT PAS UNE MAUVAISE IDÉE... ET C'EST POUR CELA QUE NOUS AVONS ÉCRIT CE LIVRE !



DANS CE QUI SUIT, NOUS ESSAIERONS DE PRÉSENTER LES ÉLÉMENTS STATISTIQUES DE LA FAÇON LA PLUS VISUELLE ET INTUITIVE POSSIBLE. TOUT CE DONT VOUS AVEZ BESOIN EST D'UN SOUPÇON DE PATIENCE, D'UN PEU DE RÉFLEXION ET D'UNE CERTAINE TOLÉRANCE À L'ALGÈBRE OU, SI CE N'EST PAS LE CAS, CE SERAIT PEUT-ÊTRE UN BON PRÉALABLE À CE COURS !

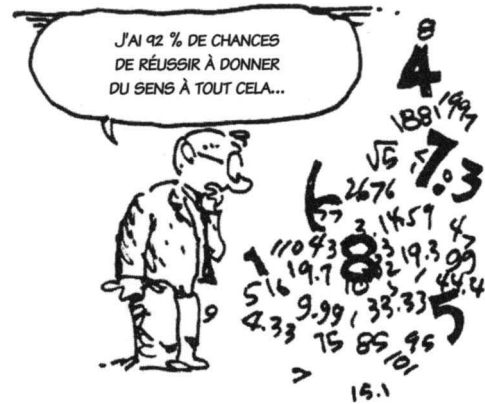


Chapitre 2

Statistiques descriptives



LES **DONNÉES** CONSTITUENT LA MATIÈRE BRUTE DES STATISTICIENS. NOUS UTILISONS DES NOMBRES POUR TRANSCRIRE LA RÉALITÉ. TOUS LES PROBLÈMES STATISTIQUES IMPLIQUENT UNE COLLECTE, UNE DESCRIPTION ET UNE ANALYSE DES DONNÉES **OU** UNE RÉFLEXION SUR CETTE COLLECTE, SUR CETTE DESCRIPTION ET SUR L'ANALYSE DES DONNÉES.



CE CHAPITRE SE CONCENTRE SUR LA PARTIE **DESCRIPTIVE** DES DONNÉES. COMMENT REPRÉSENTER LES DONNÉES DE FAÇON PRATIQUE ET UTILE? COMMENT SOULIGNER DES TENDANCES À PARTIR D'UNE ACCUMULATION DE DONNÉES BRUTES? COMMENT SYNTHÉTISER ET RÉSUMER LA FORME BASIQUE DE CES DONNÉES?



EH BIEN, POUR DÉCRIRE LES DONNÉES, NOUS AVONS BESOIN AVANT TOUT DE DISPOSER DE DONNÉES... ALORS, COMMENÇONS LEUR COLLECTE!



VOICI DES DONNÉES RÉELLEMENT COLLECTÉES LORS D'UNE EXPÉRIENCE EN CLASSE. ELLES CONCERNENT LE POIDS EN LIVRES DE 92 ÉTUDIANTS DE L'UNIVERSITÉ DE PENN STATE :



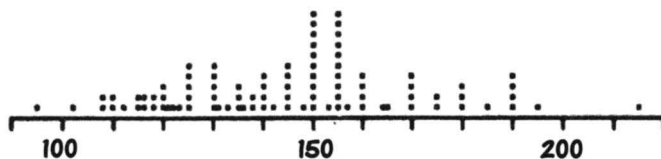
POIDS DES 57 ÉTUDIANTS

140 145 160 190 155 165 150 190 195 138 160 155 153 145 170 175 175 170 180
135 170 157 130 185 190 155 170 155 215 150 145 155 155 150 155 150 180 160
135 160 130 155 150 148 155 150 140 180 190 145 150 164 140 142 136 123 155

POIDS DES 35 ÉTUDIANTES

140 120 130 138 121 125 116 145 150 112 125 130 120 130 131 120 118 125 135
125 118 122 115 102 115 150 110 116 108 95 125 133 110 150 108

REVENONS-EN AUX CHOSES SÉRIEUSES ET TRAÇONS UN GRAPHIQUE DE POINTS : POUR CHAQUE ÉTUDIANT ON ASSOCIE UN POINT CORRESPONDANT À SON POIDS EN LIVRES. ET ON EMPILE LES POINTS LORSQUE PLUSIEURS ÉTUDIANTS ONT DONNÉ LA MÊME RÉPONSE.



POIDS EN LIVRES (1 LIVRE = 0,453 kg)



VOUS POUVEZ NOTER ICI UN **PROBLÈME** : LES AMAS DE POINTS SUR 150 ET 155 LIVRES. LES ÉTUDIANTS ONT EU TENDANCE À COMMUNIQUER LEUR POIDS EN **ARRONDISSANT À 5 LIVRES** PRÈS. DANS LA RÉALITÉ, CE TYPE D'ARRONDI PEUT NUIRE À L'ANALYSE DE TENDANCE SUR LES DONNÉES... MAIS POUR LE MOMENT, NOUS ALLONS TRAVAILLER AVEC CELLES-CI.

ON PEUT RÉSUMER LES DONNÉES EN UTILISANT UN **TABEAU DES EFFECTIFS**.
 ON REGROUPE ALORS LES DONNÉES PAR TRANCHES (NOMMÉES « CLASSES »)
 ET ON COMPTE LE NOMBRE DE POIDS D'ÉTUDIANTS DANS CHAQUE INTERVALLE.
 L'EFFECTIF CORRESPOND AU NOMBRE D'ÉTUDIANTS POUR CHAQUE CLASSE.
 LA FRÉQUENCE EST ALORS LA PROPORTION DES POIDS DANS CHAQUE INTERVALLE.
 IL S'AGIT DES EFFECTIFS DIVISÉS PAR LE NOMBRE TOTAL D'ÉTUDIANTS.

CLASSE	CENTRE DE CLASSE	EFFECTIF	FRÉQUENCE
87,5-102,4	95	2	0,022
102,5-117,4	110	9	0,098
117,5-132,4	125	19	0,206
132,5-147,4	140	17	0,185
147,5-162,4	155	27	0,293
162,5-177,4	170	8	0,087
177,5-192,4	185	8	0,087
192,5-207,4	200	1	0,011
207,5-222,4	215	1	0,011
TOTAL		92	1,000

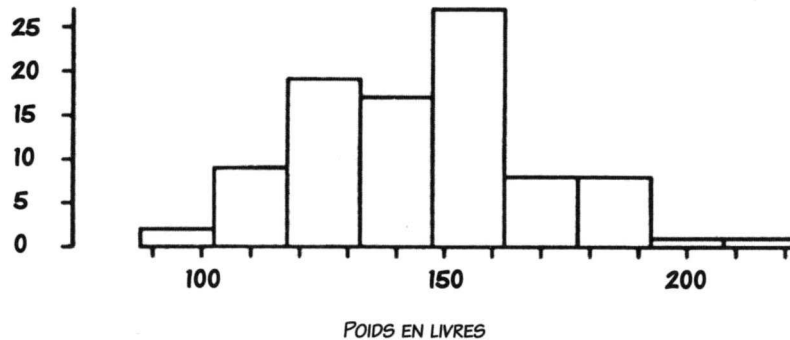
NOTE : NOUS AVONS FIXÉ, POUR LES INTERVALLES, DES BORNES ÉLOIGNÉES
 DES PROBLÈMES D'ARRONDIS À 5 LIVRES PRÈS. CELA PERMET DE LIMITER
 LES BIAIS VENANT DES DÉCLARATIONS DES ÉTUDIANTS.

DIRECTIVES POUR CRÉER DES CLASSES :

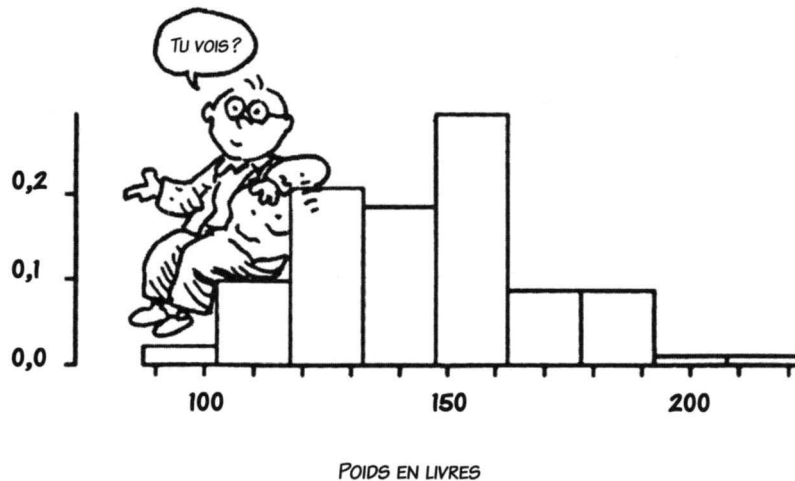
- 1)** UTILISER DES CLASSES DE MÊME AMPLITUDE AVEC DES CENTRES AUX VALEURS ARRONDIES APPROPRIÉES.
- 2)** POUR UN PETIT NOMBRE DE DONNÉES, UTILISER UN PETIT NOMBRE DE CLASSES.
- 3)** POUR UN GRAND NOMBRE DE DONNÉES, UTILISER PLUS DE CLASSES.



DANS LE TABLEAU DES EFFECTIFS, NOUS MONTRONS COMBIEN DE DONNÉES SE TROUVENT « AUTOUR » DE CHAQUE CENTRE DE CLASSE. ON PEUT AUSSI FAIRE UN DESSIN POUR REPRÉSENTER CE TYPE D'INFORMATION. LE DIAGRAMME EN BARRE ASSOCIÉ EST APPELÉ **HISTOGRAMME**. CHAQUE BARRE REPRÉSENTE UNE CLASSE ET EST CENTRÉE SUR LE CENTRE DE CLASSE. LA HAUTEUR DU RECTANGLE EST LE NOMBRE DE DONNÉES DANS LA CLASSE.



ON PEUT AUSSI DESSINER UN **HISTOGRAMME DES FRÉQUENCES**, EN REPORTANT LA FRÉQUENCE POUR CHAQUE CLASSE DE POIDS. LE GRAPHE EST EXACTEMENT LE MÊME, SEULE L'ÉCHELLE DE L'AXE VERTICAL CHANGE.



LE STATISTICIEN JOHN TUKEY (1915-2000)
A INVENTÉ UNE MÉTHODE RAPIDE
POUR RÉSUMER DES DONNÉES TOUT
EN CONSERVANT LES INFORMATIONS
BRUTES. IL S'AGIT DU DIAGRAMME
BRANCHE ET FEUILLE.



POUR NOS DONNÉES DE POIDS, LA BRANCHE
EST UNE COLONNE DE NOMBRES,
CORRESPONDANT AUX POIDS EN LIVRES
RANGÉS PAR DIZAINES (C'EST-À-DIRE
EN OMETTANT LE DERNIER CHIFFRE).

9
10
11
12
13
14
15
16
17
18
19
20
21

SOIT 90 LIVRES, 100 LIVRES,
ETC.



POUR LE POIDS DES ÉTUDIANTES
DE LA PAGE 9, ON AJOUTE LE DERNIER
CHIFFRE DE CHAQUE POIDS DANS
LA LIGNE APPROPRIÉE.

BRANCHE : FEUILLES

9 : 5
10 : 288
11 : 628855060
12 : 0155005525
13 : 0800153
14 : 05
15 : 000
16 :
17 :
18 :
19 :
20 :
21 :

CELA SIGNIFIE
QU'IL Y A
DES POIDS DE 95,
102, 108, 108,
ETC.



ENSUITE, ON PEUT RANGER LES « FEUILLES »
PAR ORDRE CROISSANT.

BRANCHE : FEUILLES

9 : 5
10 : 288
11 : 002556688
12 : 000125555
13 : 0001358
14 : 05
15 : 000
16 :
17 :
18 :
19 :
20 :
21 :

DE LA MÊME MANIÈRE, ON PEUT FAIRE
LE DIAGRAMME BRANCHE ET FEUILLE
DES POIDS DES ÉTUDIANTS HOMMES.

9 :
10 :
11 :
12 : 13
13 : 005568
14 : 000255558
15 : 000000035555555557
16 : 000045
17 : 000055
18 : 0005
19 : 00005
20 :
21 : 5

TOUS CES 0 ET CES 5
INDIQUENT CLAIREMENT LE BIAS DÙ AUX
INDICATIONS ARRONDIES DES ÉTUDIANTS.

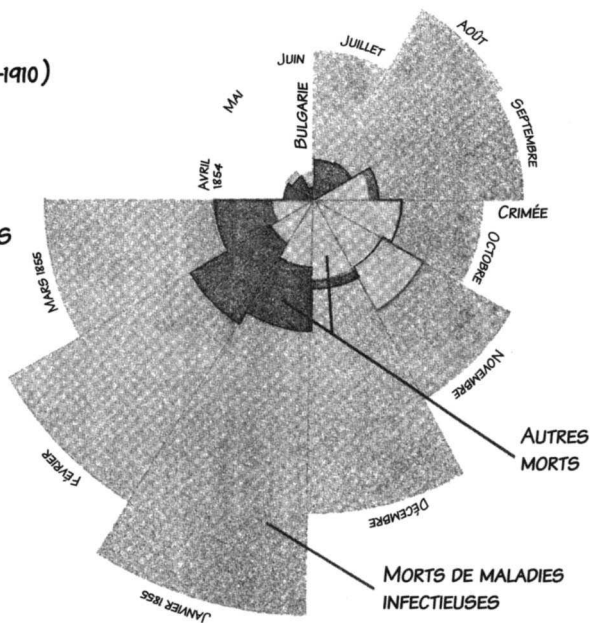


UN BON AFFICHAGE
GRAPHIQUE DOIT ÊTRE
À LA FOIS ARTISTIQUE
ET SCIENTIFIQUE.



ET PARFOIS POLITIQUE !

L'INFIRMIÈRE DÉVOUÉE
FLORENCE NIGHTINGALE (1820-1910)
A COMPILÉ LES **STATISTIQUES**
DE MORTALITÉ DANS
LES HÔPITAUX MILITAIRES
BRITANNIQUES. ELLE A PRODUIT
DES HISTOGRAMMES TROUBLANTS
COMME CELUI-CI : LE RAYON
AUGMENTE AVEC LE NOMBRE
DE SOLDATS BRITANNIQUES
MORTS LORS DE LA GUERRE
DE CRIMÉE (QUE CE SOIT
DANS LES HÔPITAUX OU SUR
LE CHAMP DE BATAILLE).
LA PLUPART DE CES MORTS
ÉTAIENT, ET DE LOIN, DUES
À DES MALADIES INFECTIEUSES.



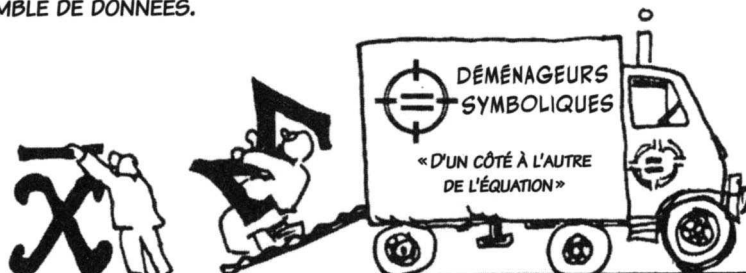
SES EFFORTS STATISTIQUES
ONT DIRECTEMENT PERMIS
D'AMÉLIORER LES CONDITIONS
À L'HÔPITAL. CE QUI
A ENTRAÎNÉ UNE RÉDUCTION
DU TAUX DE DÉCÈS.



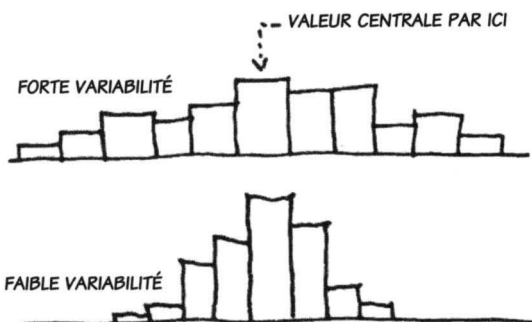
SAUVÉ PAR
LES STATISTIQUES !

RÉSUMÉ STATISTIQUE NUMÉRIQUE

MAINTENANT, PASSONS DES IMAGES AUX FORMULES. NOTRE OBJECTIF EST D'OBTENIR DES MESURES SIMPLES DES CARACTÉRISTIQUES RUDIMENTAIRES D'UN ENSEMBLE DE DONNÉES.



TOUS LES TYPES DE MESURES ONT DEUX DIMENSIONS IMPORTANTES : LA **VALEUR** OU **TENDANCE CENTRALE**, ET LA **VARIABILITÉ** AUTOUR DE CETTE VALEUR. CETTE IDÉE APPARAÎT SUR LES DEUX HISTOGRAMMES HYPOTHÉTIQUES REPRÉSENTÉS.



ON PEUT FAIRE BEAUCOUP DE CHOSES AVEC UN PETIT PEU DE NOTATION. SUPPOSONS QUE L'ON AIT UNE SÉRIE D'OBSERVATIONS... n POUR ÊTRE PRÉCIS... ALORS ON PEUT NOTER :

$$x_1, x_2, x_3 \dots x_n$$

LES DIFFÉRENTES VALEURS OBSERVÉES. AINSI, n EST LE NOMBRE TOTAL D'OBSERVATIONS, ET x_4 (PAR EXEMPLE) EST LA VALEUR DE LA QUATRIÈME OBSERVATION.

UNE **MATRICE** EST UN TABLEAU DE DONNÉES :

OBSERVATION	1	2	3	4	...	n
VALEUR OBSERVÉE	x_1	x_2	x_3	x_4	...	x_n



SUR UN PETIT ENSEMBLE DE $n = 5$ DONNÉES, ON PEUT TOUT FAIRE À LA MAIN.
PAR EXEMPLE, ON DEMANDE À 5 PERSONNES COMBIEN D'HEURES PAR SEMAINE
ILS PASSENT DEVANT LA **TÉLÉVISION...** ET VOICI LA MATRICE DES RÉSULTATS :

OBSERVATION	1	2	3	4	5
VALEUR OBSERVÉE	5	7	3	38	7

ALORS $x_1 = 5$, $x_2 = 7$, $x_3 = 3$, $x_4 = 38$, $x_5 = 7$

QUEL EST LE « CENTRE »
DE CES DONNÉES ? EN FAIT,
IL EXISTE DIFFÉRENTES RÉPONSES
À CETTE QUESTION. NOUS ALLONS
EN EXAMINER DEUX.



LA MOYENNE

LA **MOYENNE** EST REPRÉSENTÉE
PAR \bar{x} . ELLE EST OBTENUE
EN AJOUTANT TOUTES LES DONNÉES
OBSERVÉES ET EN DIVISANT
PAR LE NOMBRE D'OBSERVATIONS.

$$\bar{x} = \frac{\text{SOMME DES DONNÉES}}{n}$$

$$= \frac{x_1 + x_2 + \dots + x_n}{n}$$

POUR NOTRE EXEMPLE,

$$\bar{x} = \frac{5 + 7 + 3 + 38 + 7}{5} = \frac{60}{5}$$

= **12 HEURES**



IL EXISTE UNE NOTATION MATHÉMATIQUE PARTICULIÈRE POUR LA SOMME $x_1 + x_2 + \dots + x_n$. ON UTILISE DANS CE CAS LA LETTRE GRECQUE **SIGMA**.

Σ

AU LIEU DE $x_1 + x_2 + \dots + x_n$, ON ÉCRIT :

$$\sum_{i=1}^n x_i$$

ET ON LIT « SOMME DE x_i POUR i ÉGAL 1 À n ».

RÉPÉTEZ-LE DIX FOIS ET VOUS NE L'OUBLIEREZ JAMAIS...



SUPER! CETTE FOIS CELA RESSEMBLE VRAIMENT À UN BOUQUIN DE STATISTIQUES!



AINSI... LA **MOYENNE** D'UN ENSEMBLE DE DONNÉES x_i EST :

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad \text{OU} \quad \sum_{i=1}^n \frac{x_i}{n}$$

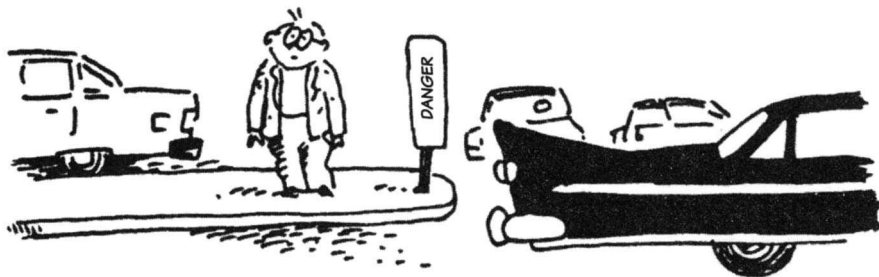
DANS LE CAS DES 92 ÉTUDIANTS DE PENN STATE, LE POIDS MOYEN EST :

$$\sum_{i=1}^{92} \frac{x_i}{92} = \frac{13\,354}{92}$$



145,15 LIVRES, SOIT ENVIRON 65,84 kg.

LA **MÉDIANE** EST UNE AUTRE MESURE DE TENDANCE CENTRALE : IL S'AGIT DU « MILIEU » QUI PARTAGE EN DEUX L'ENSEMBLE DES DONNÉES, TOUT COMME LA MÉDIANE D'UN TRIANGLE OU LE TERRE-PLEIN CENTRAL D'UNE AVENUE.



POUR TROUVER LA MÉDIANE D'UN ENSEMBLE DE DONNÉES, IL FAUT TRIER PAR ORDRE CROISSANT LES DONNÉES. LA MÉDIANE EST ALORS LA VALEUR CENTRALE.

3 5 7 7 38

LA MÉDIANE

SI LE NOMBRE D'OBSERVATIONS EST PAIR, IL N'Y A PAS VRAIMENT DE VALEUR CENTRALE. ON PREND ALORS LA MOYENNE DES DEUX VALEURS AUTOUR DU MILIEU... AINSI SI LES DONNÉES SONT :

3 5 7 7

ESPACE CENTRAL

ON FAIT LA MOYENNE ENTRE 5 ET 7 ET ON TROUVE

$$\frac{5 + 7}{2} = 6$$

CELA NOUS DONNE LA RÈGLE GÉNÉRALE DE CALCUL : ON TRIE PAR ORDRE CROISSANT LES DONNÉES.

SI LE NOMBRE D'OBSERVATIONS EST **IMPAIR**, LA MÉDIANE EST LA VALEUR CENTRALE.

SI LE NOMBRE D'OBSERVATIONS EST **PAIR**, LA MÉDIANE EST LA MOYENNE DES DEUX DONNÉES CENTRALES SITUÉES AUTOUR DU MILIEU.



POUR LES $n = 92$ POIDS D'ÉTUDIANTS,
ON PEUT CALCULER LA MÉDIANE À PARTIR
DU DIAGRAMME BRANCHE ET FEUILLE TRIÉE.
IL SUFFIT DE COMPTER JUSQU'À LA 46^e OBSERVATION.
LA MÉDIANE EST ALORS

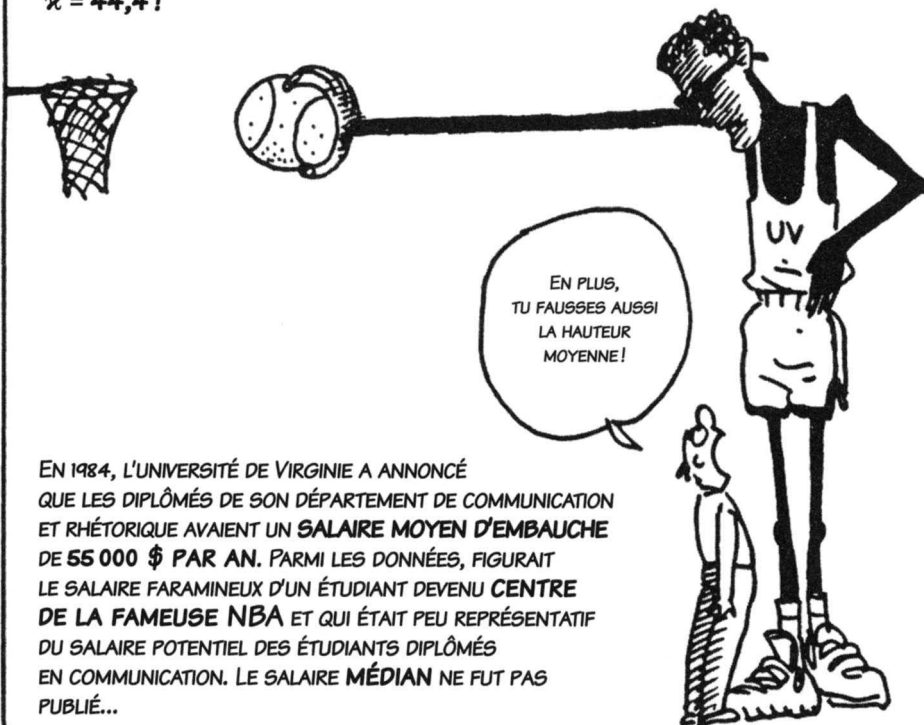
$$\frac{x_{46} + x_{47}}{2} = \frac{145 + 145}{2}$$

= **145 LIVRES** (ENVIRON 65,77 kg)

IL Y A EXACTEMENT 46 DONNÉES INFÉRIEURES
ET 46 DONNÉES SUPÉRIEURES À 145 LIVRES.

9 : 5
10 : 288
11 : 002556688
12 : 0001235555
13 : 0000013555688
14 : 00002555**55**8
15 : 000000000035555555557
16 : 000045
17 : 000055
18 : 0005
19 : 00005
20 :
21 : 5

POURQUOI UTILISONS-NOUS DEUX MESURES DIFFÉRENTES DE TENDANCE CENTRALE ?
EN FAIT, ELLES N'ONT PAS LES MÊMES PROPRIÉTÉS. PAR EXEMPLE, LA **MÉDIANE** N'EST
PAS SENSIBLE AUX DONNÉES EXTRÊMES QUI SONT ATYPIQUES PAR RAPPORT AUX AUTRES
OBSERVATIONS. SUPPOSONS QUE DANS NOTRE PETIT ÉCHANTILLON DE TÉLÉSPECTATEURS,
UNE PERSONNE REGARDE LA TÉLÉVISION 200 HEURES PAR SEMAINE. SI NOS DONNÉES
SONT $\{3, 5, 5, 7, 200\}$, LA MÉDIANE EST ENCORE 7 ALORS QUE LA MOYENNE EST MAINTENANT
 $\bar{x} = 44,4$!

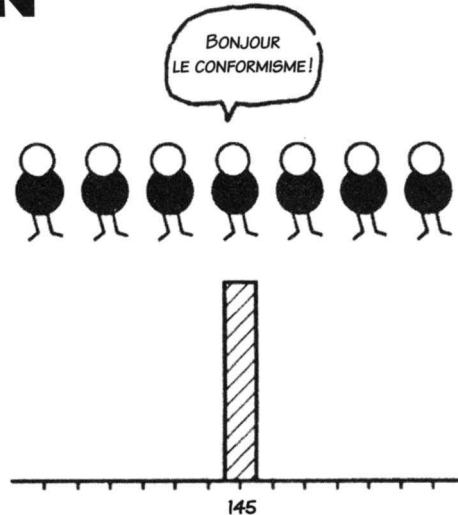


EN 1984, L'UNIVERSITÉ DE VIRGINIE A ANNONCÉ
QUE LES DIPLÔMÉS DE SON DÉPARTEMENT DE COMMUNICATION
ET RHÉTORIQUE AVAIENT UN **SALAIRE MOYEN D'EMBAUCHE**
DE **55 000 \$ PAR AN**. PARMI LES DONNÉES, FIGURAIT
LE SALAIRE FARAMINEUX D'UN ÉTUDIANT DEVENU **CENTRE**
DE LA FAMEUSE NBA ET QUI ÉTAIT PEU REPRÉSENTATIF
DU SALAIRE POTENTIEL DES ÉTUDIANTS DIPLÔMÉS
EN COMMUNICATION. LE SALAIRE **MÉDIAN** NE FUT PAS
PUBLIÉ...

Mesures de DISPERSION

(OU DE VARIABILITÉ)

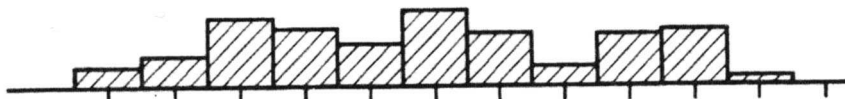
OUTRE LE CALCUL DE TENDANCE CENTRALE D'UN ENSEMBLE DE DONNÉES, NOUS SOUHAITERIONS AUSSI CONNAÎTRE LA **DISPERSION** DES DONNÉES (DE COMBIEN LES DONNÉES S'ÉLOIGNENT DU CENTRE). PAR EXEMPLE, SI LES ÉTUDIANTS PESAIENT TOUS **EXACTEMENT** 145 LIVRES, IL N'Y AURAIT AUCUNE VARIABILITÉ. NUMÉRIQUEMENT, LA VARIANCE SERAIT ÉGALE À **ZÉRO** ET L'HISTOGRAMME SERAIT TRÈS « LÉGER ».



MAIS SI BEAUCOUP D'ÉTUDIANTS ÉTAIENT TRÈS MAIGRES ET/OU D'AUTRES TRÈS LOURDS, NOUS AURIONS ÉVIDEMMENT UNE GRANDE VARIABILITÉ - PAR EXEMPLE, SI L'ÉQUIPE DE FOOTBALL AMÉRICAIN FAISAIT PARTIE DE L'ÉCHANTILLON...



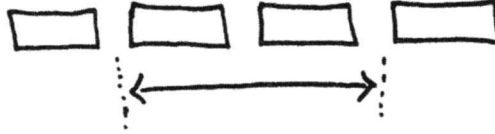
L'HISTOGRAMME SERAIT PLUS ÉTENDU, ET IL RESSEMBLERAIT À CELA :



À NOUVEAU, IL Y A PLUSIEURS FAÇONS DE DÉFINIR LA DISPERSION. L'UNE D'ELLES EST

L'ÉTENDUE INTERQUARTILE.

L'IDÉE EST DE DIVISER
LES DONNÉES EN QUATRE
GROUPE ET DE MESURER
LA DISTANCE ENTRE LES DEUX
GROUPE EXTRÊMES.



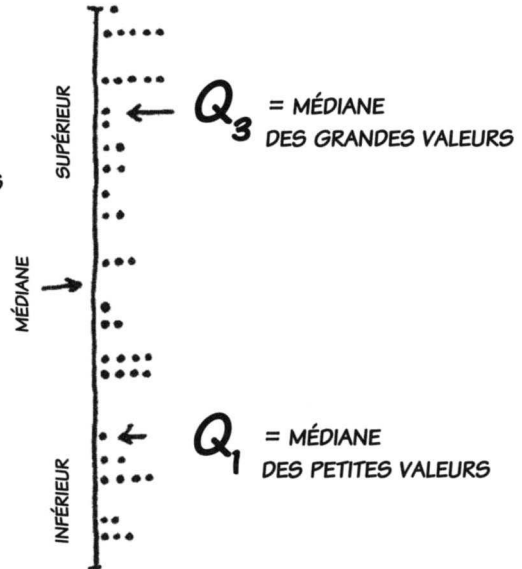
VOICI LA RECETTE :

1) ORDONNER LES DONNÉES PAR ORDRE CROISSANT.

2) DIVISER LES DONNÉES EN UN GROUPE EN DESSOUS
ET UN GROUPE AU-DESSUS DE LA MÉDIANE
(SI LA MÉDIANE CORRESPOND À UNE OBSERVATION,
LA METTRE DANS CHAQUE GROUPE).

3) CHERCHER LA MÉDIANE DU GROUPE INFÉRIEUR.
IL S'AGIT DU PREMIER QUARTILE NOTÉ Q_1 .

4) CHERCHER LA MÉDIANE DU GROUPE SUPÉRIEUR.
IL S'AGIT DU TROISIÈME QUARTILE NOTÉ Q_3 .



MAINTENANT, L'ÉTENDUE INTERQUARTILE EIQ EST LA DISTANCE,
OU LA DIFFÉRENCE, ENTRE CES DEUX QUARTILES :

$$EIQ = Q_3 - Q_1$$

VOICI LES DONNÉES DE POIDS
AVEC LES MÉDIANES DES POIDS
INFÉRIEURS ET SUPÉRIEURS
MIS EN ÉVIDENCE :

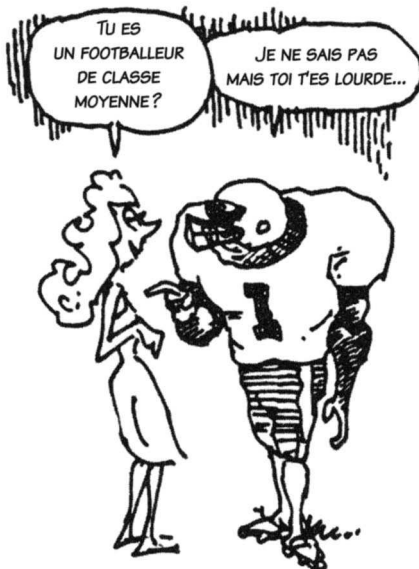
9 : 5
10 : 288
11 : 002556688
12 : 00012355555
13 : 0000013555688
14 : 000025555 58
15 : 000000000035555555557
16 : 000045
17 : 000055
18 : 0005
19 : 00005
20 :
21 : 5

Q_1
MÉDIANE
 Q_3

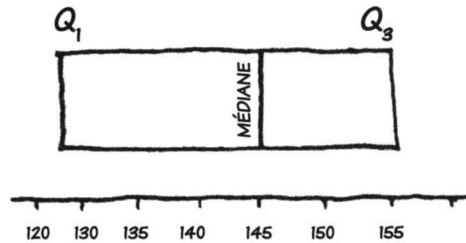
ET ON VOIT QUE

$$\begin{aligned} \text{EIQ} &= 156 - 125 \\ &= 31 \text{ LIVRES} \end{aligned}$$

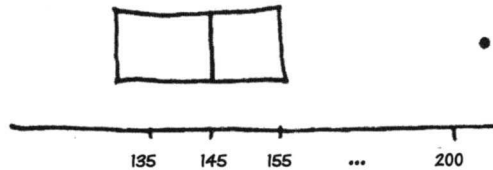
AINSI, IL S'AGIT DE LA DIFFÉRENCE ENTRE
LA MÉDIANE DES POIDS ÉLEVÉS ET CELLE
DES POIDS LÉGERS.



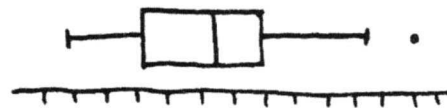
JOHN TUKEY A INVENTÉ UNE AUTRE FAÇON
DE REPRÉSENTER L'EIQ, QUE L'ON APPELLE
UNE **BOÎTE À PATTES**. LES LIMITES DE LA BOÎTE
SONT LES QUARTILES Q_1 ET Q_3 . ON TRACE
LA MÉDIANE À L'INTÉRIEUR DE LA BOÎTE.



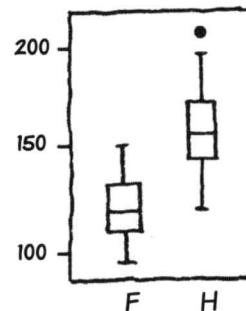
SI UN POINT EST À PLUS DE 1,5 EIQ DES LIMITES
DE LA BOÎTE, ALORS C'EST UNE **VALEUR EXTRÊME**.
ON LES REPRÉSENTE INDIVIDUELLEMENT PAR DES POINTS.



ENFIN, ON AJOUTE LES « PATTES » JUSQU'ÀUX
DERNIÈRES OBSERVATIONS QUI NE SONT PAS
DES VALEURS EXTRÊMES (C'EST-À-DIRE À MOINS
DE 1,5 EIQ).



LES BOÎTES À PATTES
SONT EXTRÊMEMENT
UTILES POUR FAIRE
RESSORTIR
LES DIFFÉRENCES
ENTRE DES GROUPES.
VOICI LES DEUX
REPRÉSENTATIONS
POUR LES ÉTUDIANTS
HOMMES ET FEMMES.



LA MESURE STANDARD DE VARIABILITÉ EST

L'ÉCART-TYPE.

CONTRAIREMENT À L'EIQ, QUI EST BASÉE SUR LES QUARTILES, L'ÉCART-TYPE MESURE LES ÉCARTS À LA **MOYENNE**. ON PEUT EN GROS RETENIR QU'IL S'AGIT DE L'ÉCART MOYEN DES DONNÉES PAR RAPPORT À LA MOYENNE \bar{x} .

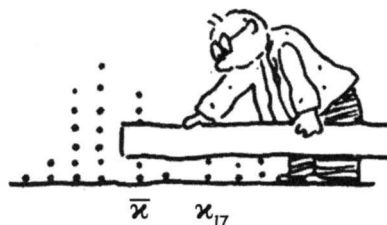


SAUF QU'EN RÉALITÉ, ON UTILISE LES **CARRÉS** DE CES ÉCARTS. AINSI, LE CARRÉ DE L'ÉCART ENTRE x_i ET \bar{x} ÉTANT $(x_i - \bar{x})^2$, ON OBTIENT :

$$\text{MOYENNE DES CARRÉS DES ÉCARTS} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

POUR DES RAISONS TECHNIQUES, ON UTILISE $n - 1$ AU DÉNOMINATEUR PLUTÔT QUE n , ET ON DÉFINIT ALORS LA **VARIANCE D'ÉCHANTILLON** s^2 COMME :

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$



POUR L'ÉCHANTILLON $\{3, 5, 7, 7, 38\}$, ON AVAIT $\bar{x} = 12$
ET $n = 5$, ON CALCULE ALORS LA VARIANCE :

$$\begin{aligned} s^2 &= \frac{(3-12)^2 + (5-12)^2 + (7-12)^2 + (7-12)^2 + (38-12)^2}{(5-1)} \\ &= \frac{81 + 49 + 25 + 25 + 676}{4} \\ &= 214 \end{aligned}$$

ICI, LA VARIANCE ÉLEVÉE REFLÈTE LA FORTE DISPERSION DES DONNÉES...



UNE MESURE DE DISPERSION DOIT ÊTRE DE MÊME UNITÉ QUE LES DONNÉES ORIGINALES. MALHEUREUSEMENT, DANS L'EXEMPLE DES POIDS, LA VARIANCE s^2 EST EN LIVRES AU CARRÉ... OUPS!



IL EST ÉVIDENT QU'IL SUFFIT DE PRENDRE LA RACINE CARRÉE, CE QUE NOUS FAISONS POUR DÉFINIR :

L'ÉCART-TYPE $s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

CE QUI DONNE POUR NOTRE ÉCHANTILLON DE DONNÉES :

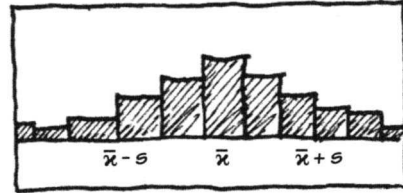
$$s = \sqrt{214} = 14,63$$



MÊME POUR UN FAIBLE NOMBRE DE DONNÉES, LES CALCULS PEUVENT ÊTRE FASTIDIEUX! HEUREUSEMENT AUJOURD'HUI, IL SUFFIT D'APPUYER SUR LE BOUTON D'UNE CALCULATRICE OU DE CONSULTER LE RAPPORT DE DONNÉES GÉNÉRÉ PAR UN LOGICIEL INCLUANT UN PACK STATISTIQUE.

Propriétés de \bar{x} et s

LA MOYENNE ET L'ÉCART-TYPE
RÉSUMENT TRÈS BIEN LES PROPRIÉTÉS
DES HISTOGRAMMES SYMÉTRIQUES
SANS VALEURS EXTRÊMES,
C'EST-À-DIRE LES HISTOGRAMMES
EN FORME DE MONTICULE.

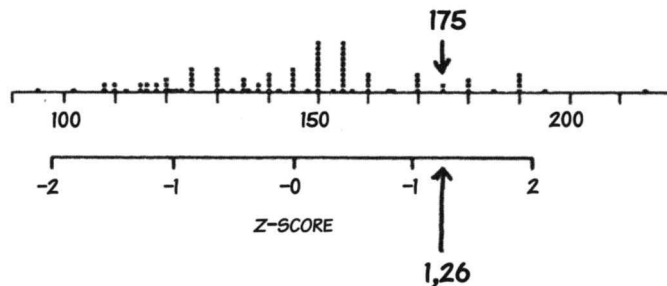


IL EST SOUVENT UTILE DE SAVOIR À **COMBIEN D'ÉCARTS-TYPES** SE TROUVE UNE DONNÉE
PAR RAPPORT À LA MOYENNE. ON DÉFINIT ALORS LE **Z-SCORE**, OU VARIABLE CENTRÉE
RÉDUITE, COMME LA DISTANCE D'UNE OBSERVATION À LA MOYENNE PAR **ÉCART-TYPE**.

$$z_i = \frac{x_i - \bar{x}}{s} \quad \text{POUR CHAQUE } i.$$



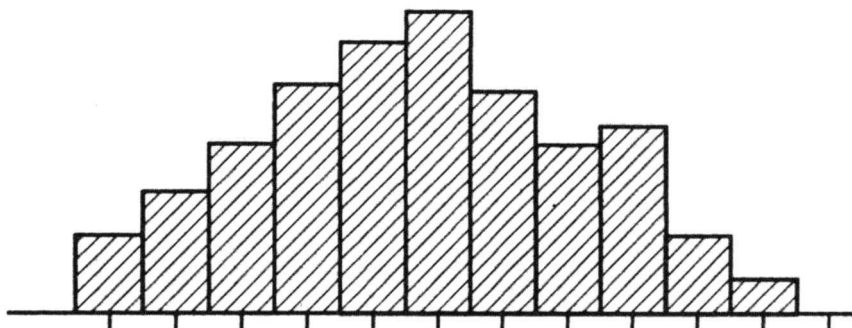
UN Z-SCORE DE + 2 SIGNIFIE QUE L'OBSERVATION EN QUESTION EST À **DEUX ÉCARTS-TYPES**
AU-DESSUS DE LA MOYENNE. POUR L'ÉCHANTILLON DE POIDS ($\bar{x} = 145,2$ ET $s = 23,7$),
ON PEUT REPRÉSENTER LES DONNÉES INITIALES SUR L'AXE DES POIDS EN LIVRES ET LEUR Z-SCORE.



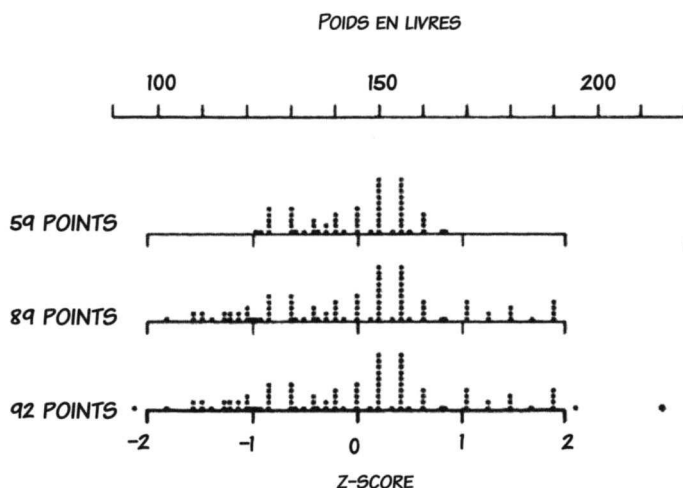
UN ÉTUDIANT QUI PÈSE 175 LIVRES A UN Z-SCORE DE $\frac{175 - 145,2}{23,7} = 1,26$

Une RÈGLE EMPIRIQUE :

POUR LA PLUPART DES ENSEMBLES DE DONNÉES DISTRIBUÉES SYMÉTRIQUEMENT EN FORME DE MONTICULE : APPROXIMATIVEMENT **68 %** DES DONNÉES SONT À UN ÉCART-TYPE DE LA MOYENNE ET **95 %** DES DONNÉES SONT À **DEUX** ÉCARTS-TYPES DE LA MOYENNE.



POUR LES POIDS DES ÉTUDIANTS, NOS DONNÉES VÉRIFIENT RELATIVEMENT LA RÈGLE : **64 %** (= 59/92) DES POIDS SONT À UN ÉCART-TYPE DE LA MOYENNE ET **97 %** DES POIDS SONT À DEUX ÉCARTS-TYPES DE LA MOYENNE.



QU'ELLE
EST MIGNONNE
CETTE VALEUR
EXTRÊME!



ET MAINTENANT,
APRÈS CES CALCULS
NUMÉRIQUES, UN PEU
DE REPOS!

ON A BEAUCOUP AVANCÉ DANS CE CHAPITRE! À PARTIR D'UN TAS DÉSORGANISÉ DE NOMBRES, NOUS AVONS :

- 1) TROUVÉ DIFFÉRENTES FAÇONS DE LES REPRÉSENTER VISUELLEMENT;
- 2) EXAMINÉ DEUX CONCEPTS DIFFÉRENTS DE TENDANCE CENTRALE : LA MÉDIANE ET LA MOYENNE;
- 3) MESURÉ LA DISPERSION DES DONNÉES AUTOUR DU CENTRE DE DEUX FAÇONS DIFFÉRENTES;
- 4) RENCONTRÉ DES HISTOGRAMMES EN FORME DE MONTICULE, ET DÉFINI Z, UNE VARIABLE QUI INDIQUE À COMBIEN D'ÉCARTS-TYPES ON SE TROUVE PAR RAPPORT À LA MOYENNE.



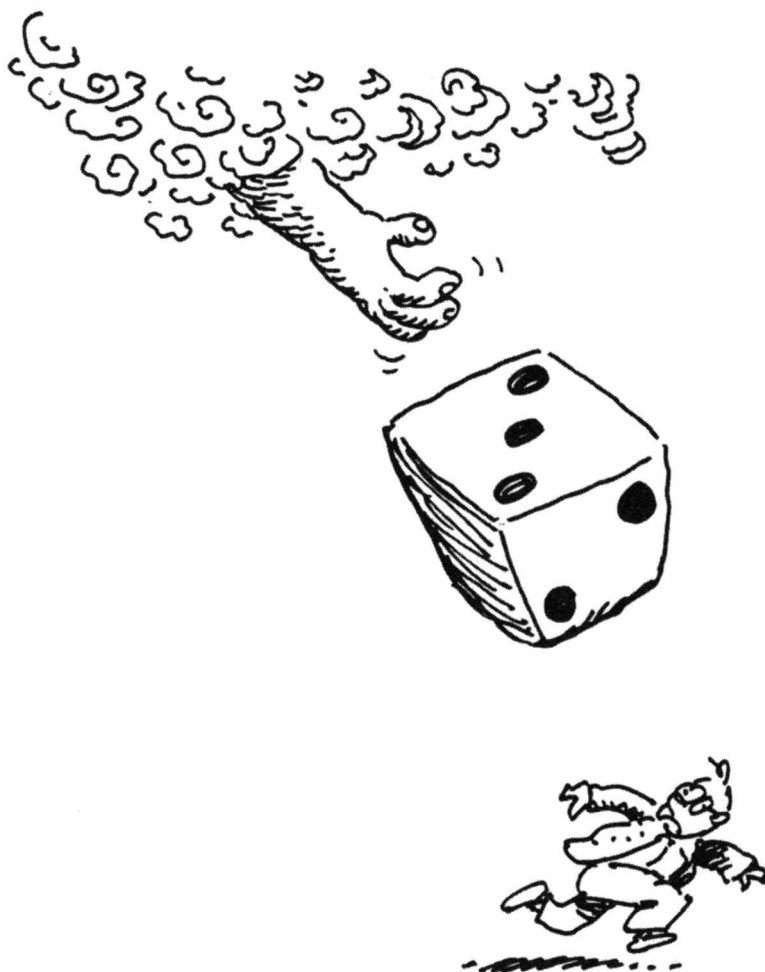
MAINTENANT, AFIN D'EXPLORER LE COMPORTEMENT DES DONNÉES DE FAÇON PLUS PRÉCISE, NOUS ALLONS FAIRE UN PETIT DÉTOUR DANS LE DOMAINE DE L'ALÉATOIRE... UN MONDE OÙ TOUT FONCTIONNE **TOUJOURS** À MERVEILLE SUR LE LONG TERME, ET OÙ LA SEULE LOI QUI TIENNE EST LA **LOI DU CASINO**.



Chapitre 3

LES PROBABILITÉS

RIEN N'EST SÛR ET CERTAIN DANS LA VIE. DANS TOUT CE QUE NOUS FAISONS, NOUS ÉVALUONS LES CHANCES DE SUCCÈS OU D'ÉCHEC, DU BUSINESS À LA MÉDECINE EN PASSANT PAR LA MÉTÉOROLOGIE. MAIS DEPUIS LES DÉBUTS DE L'AVENTURE HUMAINE, LE CALCUL DES **PROBABILITÉS**, QUI CONSISTE À ÉTUDIER FORMELLEMENT LES LOIS DU HASARD, AVANT TOUT, A ÉTÉ UTILISÉ DANS LE DOMAINE DES **JEUX D'ARGENT**.



PERSONNE NE SAIT QUAND
LES JEUX D'ARGENT SONT APPARUS
DANS L'HISTOIRE. CELA REMONTE
AU MOINS AU TEMPS DE L'ÉGYPTE
ANCIENNE OÙ DES HOMMES
ET DES FEMMES SE DIVERTISSAIENT
AVEC DES ASTRAGALES
(OS À QUATRE FACES PRÉLEVÉS
SUR LES PIEDS DES ANIMAUX).

ENTERRE-MOI
AVEC MES
ASTRAGALES...
JE VEUX DÉFIER
LA MORT!



L'EMPEREUR ROMAIN **CLAUDE 1^{er}** (10 AV. J.-C. – 54 APR. J.-C.) A ÉCRIT LE PREMIER TRAITÉ
SUR LES JEUX DE HASARD. MALHEUREUSEMENT SON LIVRE *COMMENT GAGNER AUX DÉS*
A ÉTÉ PERDU.

RÈGLE 1 : LAISSEZ
CÉSAR GAGNER
IV FOIS SUR V!



LES DÉS MODERNES ONT ÉTÉ POPULAIRES AU MOYEN ÂGE PUIS À LA RENAISSANCE.
AU XVII^e SIÈCLE, LE **CHEVALIER DE MÉRÉ** (1607-1684) FORMULA UNE ÉNIGME MATHÉMATIQUE :

QUE VAUT-IL MIEUX PARIER :
FAIRE UN SIX
EN QUATRE LANCERS
D'UN DÉ OU FAIRE
UN DOUBLE SIX
EN 24 LANCERS
D'UNE PAIRE DE DÉS ?



LE CHEVALIER RAISONNA SUR LE NOMBRE MOYEN DE SUCCÈS ET IL EN DÉDUISIT QUE, DANS LES DEUX CAS, CE NOMBRE ÉTAIT LE MÊME :

$$\text{PROBABILITÉ D'UN SIX} = \frac{1}{6}$$

$$\text{NOMBRE MOYEN EN 4 LANCERS} = 4\left(\frac{1}{6}\right) = \frac{2}{3}$$

$$\text{PROBABILITÉ D'UN DOUBLE SIX} = \frac{1}{36}$$

$$\text{NOMBRE MOYEN EN 24 LANCERS} = 24\left(\frac{1}{36}\right) = \frac{2}{3}$$

ALORS, POURQUOI PERDAIT-IL PLUS SOUVENT AVEC LE SECOND PARI ?



MÉRÉ POSA LA QUESTION À SON AMI, LE GÉNIAL **BLAISE PASCAL** (1623-1662).



BIEN QUE PASCAL AIT CESSÉ L'ÉTUDE DES MATHÉMATIQUES, JUGÉE COMME UN DIVERTISSEMENT D'ORDRE SEXUEL(!), IL ACCEPTA DE S'ATTAQUER AU PROBLÈME DE MÉRÉ.

PASCAL ÉCRIVIT À SON COLLÈGUE DE GÉNIE **PIERRE DE FERMAT** (1601-1665), ET APRÈS QUELQUES ÉCHANGES ÉPISTOLAIRES, ILS POSÈRENT LES BASES DE LA THÉORIE DES PROBABILITÉS SOUS SA FORME MODERNE – EXCEPTÉ BIEN SÛR CETTE VERSION EN BANDE DESSINÉE.

CHER PIERRE, QUELLE MERVEILLEUSE THÉORIE NOUS AURIONS, SI L'UN DE NOUS ÉTAIT CAPABLE DE DESSINER...



Définitions préliminaires

DE MÊME QU'UN PARIEUR JOUE, DE MÊME NOUS ALLONS JOUER AU SCIENTIFIQUE EN OBSERVANT LES RÉSULTATS :

UNE **expérience aléatoire** EST UN PROCÉDÉ QUI PERMET D'OBSERVER UN RÉSULTAT, OU UN ÉVÉNEMENT, DÉTERMINÉ PAR UN ALÉA.

LES **événements élémentaires** CORRESPONDENT À TOUS LES RÉSULTATS POSSIBLES DE L'EXPÉRIENCE ALÉATOIRE.

L'**espace échantillon** EST L'ENSEMBLE OU LA COLLECTION DE TOUS LES ÉVÉNEMENTS ÉLÉMENTAIRES.



PAR EXEMPLE, SI L'ÉVÉNEMENT EST UN LANCER DE PIÈCE, L'**EXPÉRIENCE ALÉATOIRE** CONSISTE À ENREGISTRER LE RÉSULTAT.



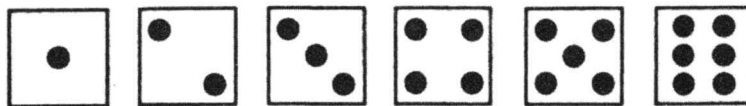
LES **ÉVÉNEMENTS ÉLÉMENTAIRES** SONT SOIT FACE, SOIT PILE.



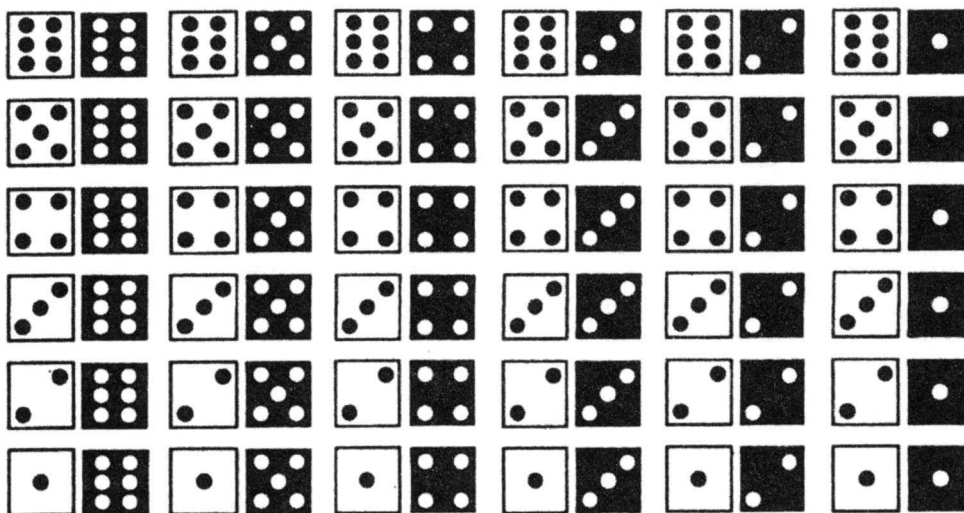
L'**ESPACE ÉCHANTILLON** EST L'ENSEMBLE SUIVANT :

$$\{F, P\}$$

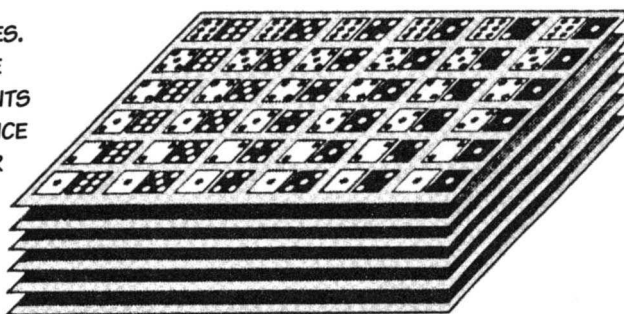

L'ESPACE ÉCHANTILLON D'UN LANCER D'UN SEUL DÉ EST UN PEU PLUS GRAND.



ET POUR LE LANCER D'UNE PAIRE DE DÉS, L'ESPACE ÉCHANTILLON RESSEMBLE À CECI (ON A PRIS UN DÉ BLANC ET UN DÉ NOIR POUR MIEUX DIFFÉRENCIER LES CAS) :

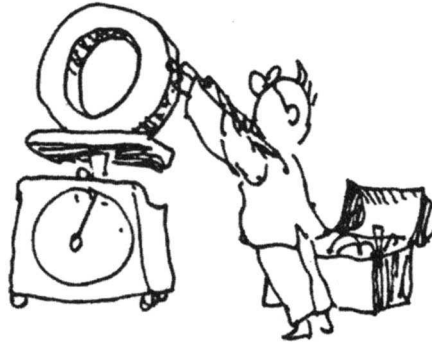


L'ESPACE ÉCHANTILLON
COMPORTE **36** (6×6)
ÉVÉNEMENTS ÉLÉMENTAIRES.
AVEC TROIS DÉS, L'ESPACE
COMPORTERAIT 216 ÉLÉMENTS
COMME DANS CETTE MATRICE
 $6 \times 6 \times 6$. ET ALORS, POUR
QUATRE DÉS ?



AU BOUT D'UN MOMENT, IL FAUT ARRÊTER
DE LISTER POUR COMMENCER À RÉFLÉCHIR...

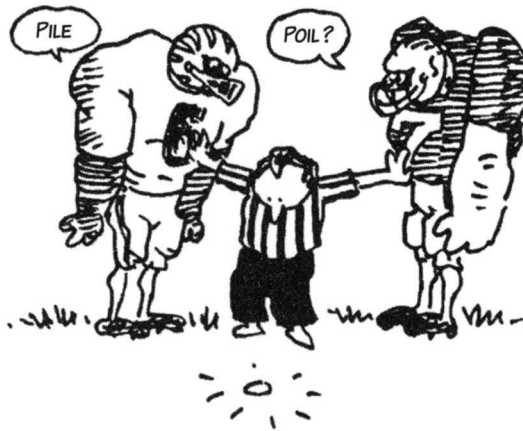
IMAGINONS MAINTENANT
UNE EXPÉRIENCE ALÉATOIRE AVEC n
RÉSULTATS ÉLÉMENTAIRES NOTÉS O_1 ,
 $O_2 \dots O_n$. ON VEUT ASSOCIER À CHAQUE
RÉSULTAT UN **POIDS NUMÉRIQUE**,
OU UNE **PROBABILITÉ** QUI MESURE
LA VRAISEMBLANCE DE L'OCCURRENCE
DE L'ÉVÉNEMENT. LA PROBABILITÉ DE O_i
SE NOTE $P(O_i)$.



PAR EXEMPLE, DANS UN LANCER
D'UNE PIÈCE (SANS TRICHER),
LES CÔTÉS FACE ET PILE SONT
ÉGALEMENT VRAISEMBLABLES.
LA PROBABILITÉ DE CHACUN
EST DE 0,5.

$$P(F) = P(P) = 0,5$$

CHAQUE ÉVÉNEMENT A LIEU
UNE FOIS SUR DEUX. DEMANDEZ
À UN JOUEUR DE FOOTBALL !

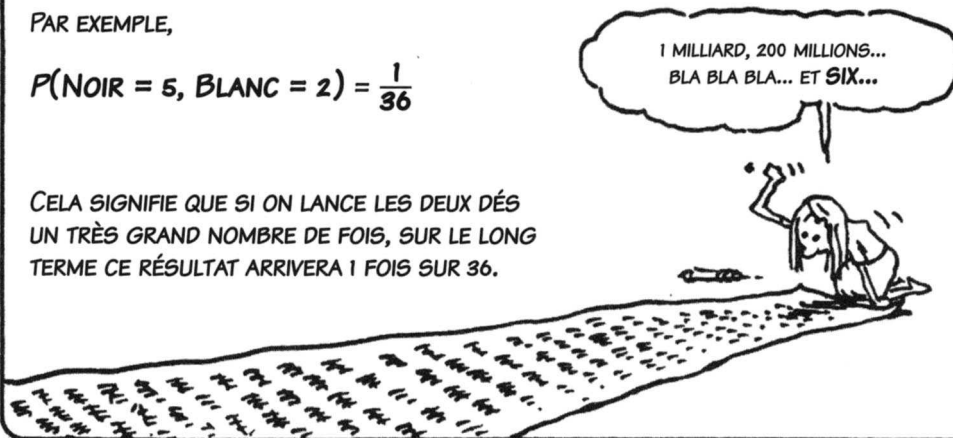


DANS UN LANCER DE **DEUX DÉS**, IL Y A **36** ÉVÉNEMENTS ÉLÉMENTAIRES.
TOUS SONT AUSSI VRAISEMBLABLES, CHAQUE PROBABILITÉ EST DONC DE $1/36$.

PAR EXEMPLE,

$$P(\text{NOIR} = 5, \text{BLANC} = 2) = \frac{1}{36}$$

CELA SIGNIFIE QUE SI ON LANCE LES DEUX DÉS
UN TRÈS GRAND NOMBRE DE FOIS, SUR LE LONG
TERME CE RÉSULTAT ARRIVERA 1 FOIS SUR 36.

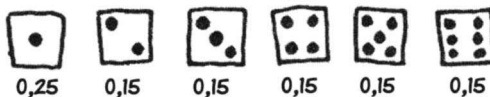


QUE SE PASSE-T-IL SI NOTRE JOUEUR
TRICHE EN UTILISANT UN DÉ PIPÉ ?
POUR SIMPLIFIER, SUPPOSONS
QUE LA FACE 1 APPARAISSE 25 %
DU TEMPS (SUR LE LONG TERME).



L'ESPACE ÉCHANTILLON
EST LE MÊME QUE CELUI
D'UN DÉ NORMAL

$\{1, 2, 3, 4, 5, 6\}$



PAR CONTRE, LES PROBABILITÉS
DIFFÈRENT. MAINTENANT $P(1) = 0,25$
ET LES PROBABILITÉS RESTANTES
DOIVENT SE SOMMER À 0,75.
SI 2, 3, 4, 5, ET 6 RESTENT ÉQUIPROBABLES,
ALORS CHAQUE FACE AURA UNE PROBABILITÉ
DE $0,15 = \frac{1}{5}$ (0,75).



EN GÉNÉRAL, LES ÉVÉNEMENTS ÉLÉMENTAIRES N'ONT PAS LA MÊME PROBABILITÉ.



MAINTENANT, QUE POUVONS-NOUS DIRE
DES PROBABILITÉS $P(O_i)$ DANS UNE EXPÉRIENCE
ALÉATOIRE QUELCONQUE. PREMIÈREMENT,

$$P(O_i) \geq 0$$

LES PROBABILITÉS NE SONT **JAMAIS
NÉGATIVES**. UNE PROBABILITÉ DE ZÉRO VEUT
DIRE QUE L'ÉVÉNEMENT N'AURA JAMAIS LIEU.
UNE VALEUR STRICTEMENT INFÉRIEURE À ZÉRO
N'AURAIT AUCUN SENS.



DEUXIÈMEMENT, S'IL EST **CERTAIN** QU'UN ÉVÉNEMENT AURA LIEU, NOUS LUI ASSIGNONS UNE PROBABILITÉ
DE 1 (SUR LE LONG TERME, C'EST LA PROPORTION DE FOIS QU'IL SE PRODUIRA).



EN PARTICULIER,
**LA PROBABILITÉ
TOTALE DE L'ESPACE
ÉCHANTILLON DOIT**
ÊTRE ÉGALE À 1. SI L'ON CONDUIT L'EXPÉRIENCE,
QUELQUE CHOSE ARRIVERA, FORCÉMENT!



EN COMBINANT LES DEUX, ON OBTIENT LES **PROPRIÉTÉS CARACTÉRISTIQUES
DES PROBABILITÉS** :

$$P(O_i) \geq 0$$

UNE PROBABILITÉ EST **POSITIVE**.

$$P(O_1) + P(O_2) + \dots + P(O_n) = 1$$

LA PROBABILITÉ TOTALE DE TOUS
LES ÉVÉNEMENTS ÉLÉMENTAIRES EST **UN**.



À LA FAÇON D'UN HABILE POLITICIEN,
NOUS AVONS ÉVITÉ CERTAINES
QUESTIONS DÉPLAISANTES,
COMME :

A) QUE **SIGNIFIE** EXACTEMENT
UNE PROBABILITÉ ?
ET B) **COMMENT** DÉFINIR
LES PROBABILITÉS DES ÉVÉNEMENTS
ÉLÉMENTAIRES ?

HEU, HEU...
DISCUTONS PLUTÔT
DE QUELQUE CHOSE
DE PLUS SIMPLE COMME
LE RÉCHAUFFEMENT
CLIMATIQUE...

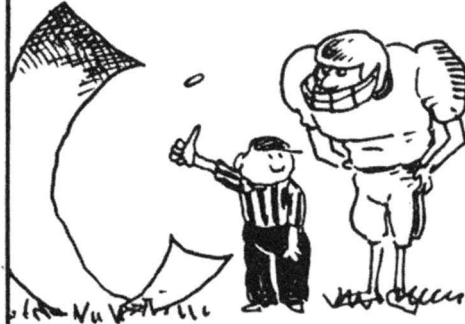


VOICI QUELQUES-UNES DES APPROCHES POSSIBLES :

LES PROBABILITÉS **classiques** :
ELLES REPOSENT SUR DES IDÉES GÉNÉRALES
SUR LES JEUX. L'HYPOTHÈSE FONDAMENTALE
EST QUE LE JEU EST ÉQUITABLE ET QUE TOUS
LES ÉVÉNEMENTS ÉLÉMENTAIRES ONT LA MÊME
PROBABILITÉ.



LES **fréquences (relatives)** : QUAND
UNE EXPÉRIENCE PEUT ÊTRE RÉPÉTÉE, ALORS
LA PROBABILITÉ D'UN ÉVÉNEMENT
EST LA PROPORTION D'OCCURRENCES
DE L'ÉVÉNEMENT SUR LE LONG TERME.



LES PROBABILITÉS **personnelles** :
LA PLUPART DES ÉVÉNEMENTS DE LA VIE
NE SONT PAS **RÉPÉTABLES**. UNE
PROBABILITÉ PERSONNELLE EST UNE ESTIMATION
INDIVIDUELLE ET PERSONNELLE
DE LA VRAISEMBLANCE D'UN ÉVÉNEMENT.
SI UN PARIEUR ESTIME QU'UN CHEVAL
A PLUS DE 50 % DE CHANCES DE GAGNER,
IL MISERA SUR CE CHEVAL.



UN **OBJECTIVISTE** UTILISERA SOIT LA DÉFINITION
CLASSIQUE, SOIT LES FRÉQUENCES COMME
PROBABILITÉS. UN **SUBJECTIVISTE**
OU UN **BAYÉSIEEN*** APPLIQUERA LES LOIS
FORMELLES DU HASARD POUR SES PROBABILITÉS
PERSONNELLES OU LES VÔTRES.



*DE THOMAS BAYES (VOIR PAGE 47).

Les opérations de base

JUSQU'À PRÉSENT, NOUS AVONS DISCUTÉ UNIQUEMENT DES PROBABILITÉS DES ÉVÉNEMENTS ÉLÉMENTAIRES. EN THÉORIE, CELA SUFFIT POUR DÉCRIRE N'IMPORTE QUELLE EXPÉRIENCE ALÉATOIRE, MAIS DANS LA PRATIQUE, C'EST PEU MANIABLE. PAR EXEMPLE, L'OCCURRENCE D'UN JET DE SOMME 7 AVEC DEUX DÉS N'EST PAS UN ÉVÉNEMENT ÉLÉMENTAIRE. NOUS INTRODUISONS DONC UNE NOUVELLE IDÉE.



UN **ÉVÉNEMENT** EST UN ENSEMBLE D'ÉVÉNEMENTS ÉLÉMENTAIRES. LA PROBABILITÉ D'UN ÉVÉNEMENT EST LA SOMME DES PROBABILITÉS DES ÉVÉNEMENTS ÉLÉMENTAIRES DE L'ENSEMBLE. VOICI DES EXEMPLES D'ÉVÉNEMENTS LORSQU'ON LANCE DEUX DÉS :

DESCRIPTION DE L'ÉVÉNEMENT	ÉVÉNEMENTS ÉLÉMENTAIRES RÉALISANT L'ÉVÉNEMENT	PROBABILITÉS
A : SOMME DES DÉS = 3	$\{(1,2), (2,1)\}$	$P(A) = \frac{2}{36}$
B : SOMME DES DÉS = 6	$\{(1,5), (2,4), (3,3), (4,2), (5,1)\}$	$P(B) = \frac{5}{36}$
C : DÉ BLANC = 1	$\{(1,1), (1,2), (1,3), (1,4), (1,5), (1,6)\}$	$P(C) = \frac{6}{36}$
D : DÉ NOIR = 1	$\{(1,1), (2,1), (3,1), (4,1), (5,1), (6,1)\}$	$P(D) = \frac{6}{36}$



ET QUAND
EST-CE QUE
JE RÉCUPÈRE
MA CHEMISE ?

L'AVANTAGE D'UTILISER DES ÉVÉNEMENTS, ET PLUS SEULEMENT DES ÉVÉNEMENTS ÉLÉMENTAIRES, EST DE POUVOIR LES **COMBINER** POUR CRÉER DE NOUVEAUX ÉVÉNEMENTS EN UTILISANT DES OPÉRATEURS LOGIQUES. LES MOTS-CLÉS UTILISÉS SONT **ET**, **OU** ET **NON**.



AINSI, ÉTANT DONNÉ DEUX ÉVÉNEMENTS E ET F , ON PEUT CRÉER LES NOUVEAUX ÉVÉNEMENTS :

E et F : OÙ LES ÉVÉNEMENTS E ET F ONT LIEU EN MÊME TEMPS.

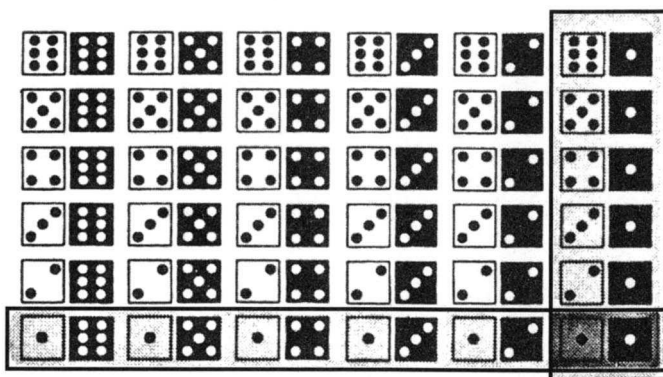
E ou F : OÙ SOIT E A LIEU, SOIT F , SOIT LES DEUX EN MÊME TEMPS.

non E : OÙ L'ÉVÉNEMENT E N'A PAS LIEU.

EN COMBINANT NOS DÉFINITIONS DE BASE SUR LES PROBABILITÉS AVEC CES OPÉRATEURS LOGIQUES, NOUS OBTENONS DE PUISSANTES FORMULES POUR MANIPULER LES PROBABILITÉS.



REVENONS À L'EXEMPLE DES LANCERS DE DÉS. SI C EST L'ÉVÉNEMENT [DÉ BLANC = 1] ET D EST L'ÉVÉNEMENT [DÉ NOIR = 1] ALORS



C ET D EST
L'INTERSECTION
(LE CHEVAUCHEMENT)
DES PARTIES
HACHURÉES
OÙ LES DEUX DÉS
VALENT 1.

C OU D EST
L'UNION DES PARTIES
HACHURÉES
OÙ AU MOINS
L'UN DES DÉS VAUT 1.

CELA ILLUSTRE LA RÈGLE D'ADDITION : POUR TOUT ÉVÉNEMENT E ET F ,

$$P(E \text{ OU } F) = P(E) + P(F) - P(E \text{ ET } F)$$

EN EFFET, LA SOMME $P(E) + P(F)$ COMPTE DEUX FOIS LES ÉVÉNEMENTS ÉLÉMENTAIRES DE E ET F , NOUS DEVONS DONC RETIRER CETTE QUANTITÉ QUI CORRESPOND À $P(E \text{ ET } F)$.

DANS L'EXEMPLE CITÉ PLUS HAUT,

$$P(C \text{ OU } D) = \frac{11}{36}$$

COMME VOUS POUVEZ LE VOIR
EN COMPTANT LES ÉVÉNEMENTS
ÉLÉMENTAIRES. DE MÊME,

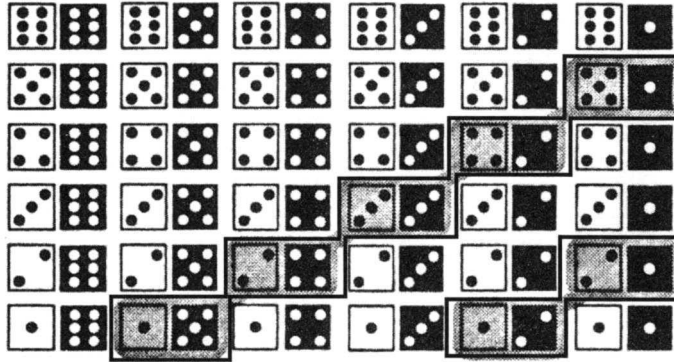
$$P(C \text{ ET } D) = \frac{1}{36}$$

ET ON CONFIRME ALORS LA FORMULE :

$$\begin{aligned} P(C) + P(D) - P(C \text{ ET } D) \\ = \frac{6}{36} + \frac{6}{36} - \frac{1}{36} &= \frac{11}{36} \\ &= P(C \text{ OU } D) \end{aligned}$$



PARFOIS LE CHEVAUCHEMENT OU L'INTERSECTION E ET F EST VIDE, ET LES DEUX ÉVÉNEMENTS NE PARTAGENT AUCUN ÉVÉNEMENT ÉLÉMENTAIRE. DANS CE CAS, ON DIT QUE E ET F SONT **MUTUELLEMENT EXCLUSIFS** DE SORTE QUE $P(E \text{ ET } F) = 0$. DANS LE DESSIN SUIVANT, ON VOIT QUE LES ÉVÉNEMENTS A [SOMME DES DÉS = 3] ET B [SOMME DES DÉS = 6] SONT MUTUELLEMENT EXCLUSIFS.



POUR DES ÉVÉNEMENTS MUTUELLEMENT EXCLUSIFS, ON A UNE LOI SPÉCIALE D'ADDITION : SI E ET F SONT MUTUELLEMENT EXCLUSIFS ALORS

$$P(E \text{ OU } F) = P(E) + P(F)$$

$$\text{ET ON PEUT VÉRIFIER } P(A \text{ OU } B) = \frac{7}{36} = \frac{2}{36} + \frac{5}{36} = P(A) + P(B)$$

ET ENFIN UNE RÈGLE DE SOUSTRACTION : POUR TOUT ÉVÉNEMENT E ,

$$P(E) = 1 - P(\text{NON } E)$$

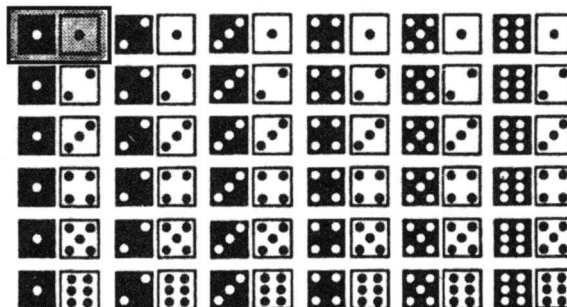
CETTE RÈGLE EST TRÈS UTILE LORSQUE $P(\text{NON } E)$ EST PLUS FACILE À CALCULER QUE $P(E)$. PAR EXEMPLE, SOIT E L'ÉVÉNEMENT [PAS DE DOUBLE 1]. L'ÉVÉNEMENT $\text{NON } E$, [RÉUSSIR UN DOUBLE 1], A UNE PROBABILITÉ DE $P(\text{NON } E) = 1/36$.

AINSI

$$P(E) = 1 - P(\text{NON } E)$$

$$= 1 - \frac{1}{36}$$

$$= \frac{35}{36}$$



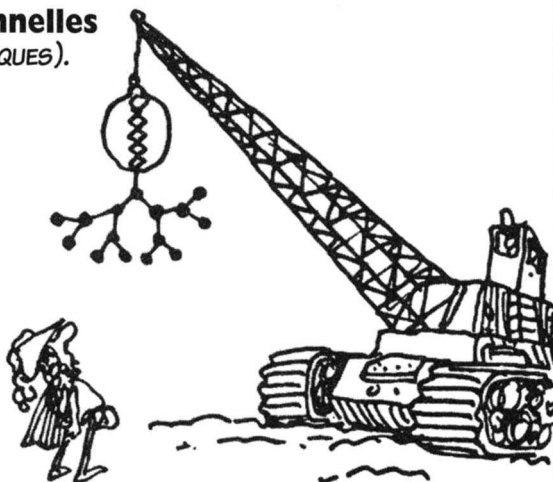
SERAIT-IL
MAINTENANT
POSSIBLE
DE RÉSOUDRE
MON PROBLÈME?
IL FAIT
FROID...



LES FORMULES QUE NOUS AVONS
PRÉSENTÉES SONT, EN FAIT, ADAPTÉES
POUR RÉSOUDRE LE PROBLÈME DE MÉRÉ.
MAIS CE N'EST PAS SIMPLE!
(ON POURRAIT ESSAYER DE RÉPONDRE
À UNE QUESTION PLUS FACILE : QUELLE EST
LA PROBABILITÉ DE LANCER UN 6 SUR DEUX
LANCERS DE DÉS ?) NOUS AVONS BESOIN
DE PLUS D'OUTILS !

NOUS INTRODUISONS DONC
LES **probabilités conditionnelles**
(UN CONCEPT ESSENTIEL EN STATISTIQUES).

WAOUH!
ÇA A L'AIR
D'ÊTRE
DU LOURD !



SUPPOSONS QUE NOUS MODIFIONS NOTRE EXPÉRIENCE EN LANÇANT LE DÉ BLANC
AVANT LE DÉ NOIR. QUELLE EST LA PROBABILITÉ DE L'ÉVÉNEMENT **A**, C'EST-À-DIRE
QUE LA SOMME DES FACES SOIT 3 ?

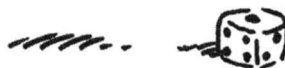


AVANT LE LANCEMENT
DES DÉS, CETTE
PROBABILITÉ EST

$$P(A) = \frac{2}{36}$$



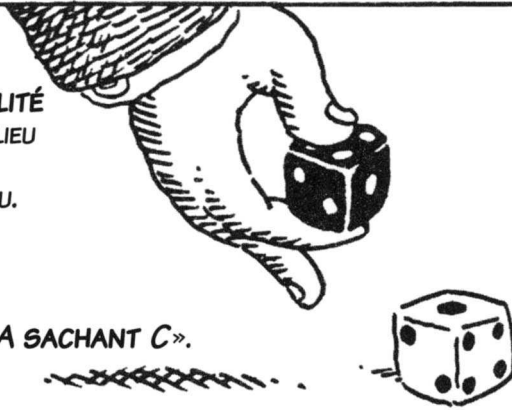
SUPPOSONS MAINTENANT
QUE LE DÉ BLANC TOMBE
SUR LE 1 (ÉVÉNEMENT **C**).
QUELLE EST MAINTENANT
LA PROBABILITÉ DE **A** ?



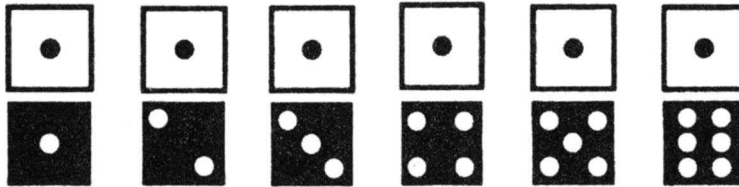
ON APPELLE CELA LA **PROBABILITÉ CONDITIONNELLE** QUE **A** AIT LIEU SOUS LA **CONDITION** QUE L'ÉVÉNEMENT **C** AIT DÉJÀ EU LIEU. ON L'ÉCRIT

$$P(A|C)$$

ET ON DIT « LA PROBABILITÉ DE **A** SACHANT **C** ».



AVANT DE JETER UN DÉ, L'ESPACE ÉCHANTILLON CONTIENT 36 ÉLÉMENTS. MAIS UNE FOIS QUE L'ÉVÉNEMENT **C** A EU LIEU, LE RÉSULTAT APPARTIENDRA À L'ESPACE ÉCHANTILLON RÉDUIT VÉRIFIANT **C**.



DANS L'ESPACE ÉCHANTILLON RÉDUIT DE SIX ÉVÉNEMENTS ÉLÉMENTAIRES, SEUL UN RÉSULTAT (1,2) DONNE UNE SOMME À 3. AINSI LA PROBABILITÉ CONDITIONNELLE EST DE 1/6.

TU VOIS COMME LES PROBABILITÉS CHANGENT LORSQUE LE MONDE ÉVOLUE ?



EN GÉNÉRAL, POUR DÉTERMINER LA PROBABILITÉ CONDITIONNELLE $P(E|F)$, ON EXAMINE L'ÉVÉNEMENT **E** ET **F** EN TANT QUE PARTIE DE L'ESPACE ÉCHANTILLON RÉDUIT VÉRIFIANT **F**.



NOUS ALLONS TRANSCRIRE
CELA EN DÉFINITION
FORMELLE : LA **PROBABILITÉ
CONDITIONNELLE
DE E SACHANT F** EST :

$$P(E|F) = \frac{P(E \text{ ET } F)}{P(F)}$$

LA FORMULE VÉRIFIE DEUX RÉSULTATS INTUITIFS :

$$P(E|E) = 1 \quad (\text{SI } E \text{ A EU LIEU, ALORS } E \text{ EST CERTAIN}).$$

LORSQUE E ET F SONT MUTUELLEMENT
EXCLUSIFS

$$P(E|F) = 0 \quad (\text{SI } F \text{ A EU LIEU, ALORS } E \text{ EST IMPOSSIBLE}).$$



AVEC LES DÉS, C'EST

$$\frac{P(A \text{ ET } C)}{P(C)} = \frac{\frac{1}{36}}{\frac{1}{6}} = \frac{1}{6}$$

EN RÉARRANGEANT LES TERMES, NOUS OBTENONS UNE **LOI MULTIPLICATIVE** :

$$P(E \text{ ET } F) = P(E|F)P(F)$$

CE QUE NOUS AIMERIONS SIMPLIFIER EN UNE LOI MULTIPLICATIVE « SPÉCIALE »
SOUS LA CONDITION QUE $P(E|F) = P(E)$. CE SERAIT EXCELLENT !



ET AVANT QUE VOUS
NE VOUS PRÉCIPITIEZ
SUR LA PROCHAINE PAGE, NOTEZ
QU'EN INTERVERTISSANT E ET F
NOUS OBTENONS QUE
 $P(F)P(E|F) = P(E)P(F|E)$.

L'INDÉPENDANCE et la loi multiplicative spéciale

DEUX ÉVÉNEMENTS E ET F SONT DITS **INDÉPENDANTS** L'UN DE L'AUTRE SI L'OCCURRENCE DE L'UN N'A **AUCUNE INFLUENCE** SUR L'OCCURRENCE DE L'AUTRE. PAR EXEMPLE, LE RÉSULTAT DU LANCER D'UN DÉ N'A AUCUNE INCIDENCE SUR LE RÉSULTAT DE L'AUTRE DÉ (À MOINS QU'ILS NE SOIENT COLLÉS ENTRE EUX, OU MAGNÉTIQUEMENT RELIÉS, ETC.!).



EN TERMES DE PROBABILITÉ CONDITIONNELLE, CELA IMPLIQUE QUE $P(E) = P(E|F)$ OU DE FAÇON ÉQUIVALENTE $P(F) = P(F|E)$. LORSQUE E ET F SONT INDÉPENDANTS, NOUS OBTENONS LA **LOI MULTIPLICATIVE SPÉCIALE** :

$$P(E \text{ ET } F) = P(E)P(F)$$

VÉRIFIONS MAINTENANT L'INDÉPENDANCE DES DÉS EN UTILISANT LES FORMULES. C EST L'ÉVÉNEMENT [DÉ BLANC = 1] ET D EST L'ÉVÉNEMENT [DÉ NOIR = 1]. NOUS AVONS ALORS

$$P(C|D) = \frac{P(C \text{ ET } D)}{P(D)} = \frac{\frac{1}{36}}{\frac{1}{6}} = \frac{1}{6} = P(C)$$

PAR CONTRE, SI LE DÉ BLANC VAUT 1, CELA AFFECTE **ÉVIDEMMENT** LA PROBABILITÉ QUE LA SOMME DES DÉS DONNE 3!

$$P(A|C) = \frac{P(A \text{ ET } C)}{P(C)} = \frac{P(1,2)}{P(C)} = \frac{\frac{1}{36}}{\frac{1}{6}} = \frac{1}{6} \neq P(A) = \frac{1}{18}$$

CES DEUX ÉVÉNEMENTS A ET C NE SONT DONC **PAS INDÉPENDANTS**.

AVANT D'ALLER PLUS LOIN, VOICI LA LISTE DES LOIS QUE NOUS AVONS OBTENUES :

LOI ADDITIVE :

$$P(E \text{ OU } F) = P(E) + P(F) - P(E \text{ ET } F)$$

LOI SPÉCIALE D'ADDITION, SI E ET F SONT MUTUELLEMENT EXCLUSIFS :

$$P(E \text{ OU } F) = P(E) + P(F)$$

LOI DE SOUSTRACTION :

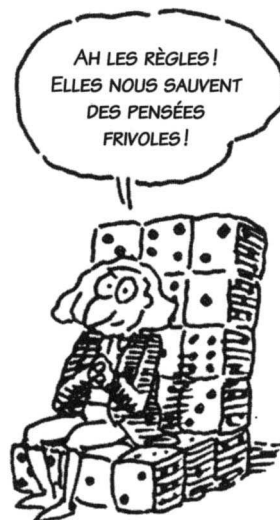
$$P(E) = 1 - P(\text{NON } E)$$

LOI MULTIPLICATIVE :

$$P(E \text{ ET } F) = P(E|F)P(F)$$

LOI SPÉCIALE DE MULTIPLICATION, SI E ET F SONT INDÉPENDANTS :

$$P(E \text{ ET } F) = P(E)P(F)$$



ET ENFIN LE PROBLÈME DE MÉRÉ... SOIT E L'ÉVÉNEMENT [OBTENIR AU MOINS UN SIX SUR QUATRE LANCERS DE DÉS], QUE VAUT $P(E)$? IL S'AGIT D'UN ÉVÉNEMENT DONT L'ÉVÉNEMENT CONTRAIRE EST PLUS SIMPLE À CALCULER : $\text{NON } E$ EST L'ÉVÉNEMENT [N'OBTENIR AUCUN SIX SUR QUATRE LANCERS].



SI A_i EST L'ÉVÉNEMENT [NE PAS OBTENIR DE SIX AU i^{e} LANCER], NOUS SAVONS QUE $P(A_i) = 5/6$. NOUS SAVONS AUSSI QUE LES LANCERS SONT INDÉPENDANTS. AINSI,

$$P(\text{NON } E) = P(A_1 \text{ ET } A_2 \text{ ET } A_3 \text{ ET } A_4)$$

LOI MULTIPLICATIVE

$$\rightarrow = \left(\frac{5}{6}\right)^4 = 0,482$$

DONC

$$P(E) = 1 - P(\text{NON } E) = 0,518$$

MAINTENANT LE SECOND PARI : SOIT F L'ÉVÉNEMENT [OBTENIR UN DOUBLE 6 SUR 24 LANCERS DE DEUX DÉS]. À NOUVEAU, **NON F** EST PLUS SIMPLE À UTILISER CAR C'EST L'ÉVÉNEMENT [N'OBTENIR **AUCUN DOUBLE SIX**].



SI B_i EST L'ÉVÉNEMENT [PAS DE DOUBLE SIX AU i^{e} LANCER] ALORS

NON F = B_1 ET B_2 ET ... B_{24} .

LA PROBABILITÉ DE CHAQUE B_i EST :

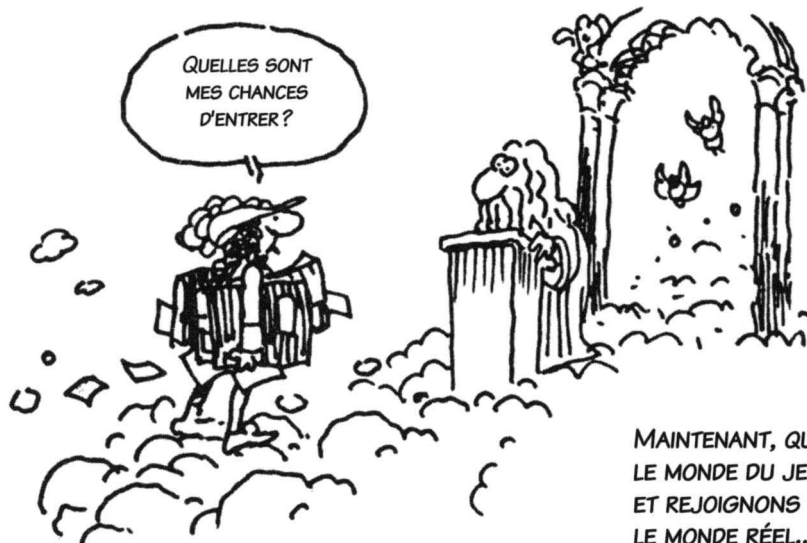
$P(B_i) = \left(\frac{35}{36}\right)$, DONC

$P(\text{NON } F) = \left(\frac{35}{36}\right)^{24} = 0,509$

PAR LA RÈGLE DE MULTIPLICATION, ON CONCLUT QUE :

$P(F) = 1 - P(\text{NON } F)$
 $= 1 - 0,509 = 0,491$

MÉRÉ AVAIT DIT À PASCAL QU'IL AVAIT OBSERVÉ QUE L'ÉVÉNEMENT F SE PRODUISAIT MOINS SOUVENT QUE L'ÉVÉNEMENT E , MAIS IL RESTAIT PERPLEXE SUR L'EXPLICATION À DONNER... NOUS EN CONCLUONS QUE MÉRÉ DEVAIT SOUVENT FAIRE LE TEST ET EN CONSERVER SOIGNEUSEMENT LES RÉSULTATS.



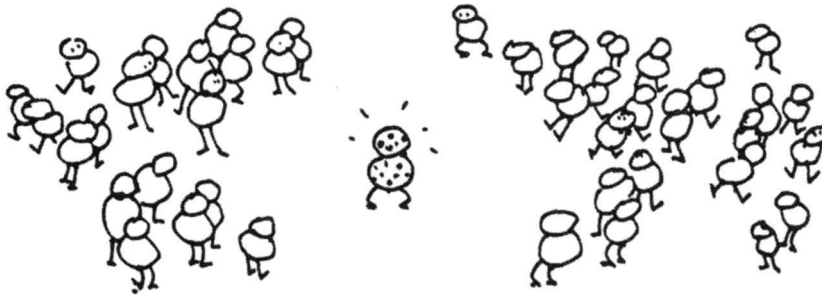
MAINTENANT, QUITTONS LE MONDE DU JEU ET REJOIGNONS LE MONDE RÉEL...

LE THÉORÈME DE BAYES et le cas des faux positifs

POUR UNE APPLICATION PLUS SÉRIEUSE
DES PROBABILITÉS CONDITIONNELLES, EXAMINONS
UN DOMAINE CONCERNANT UNE QUESTION DE VIE
OU DE MORT.



SUPPOSONS QU'UNE MALADIE MORTELLE INFECTE UN INDIVIDU SUR 1000
DANS UNE POPULATION DONNÉE...



ET SUPPOSONS QU'IL EXISTE UN BON TEST, MAIS IMPARFAIT, POUR DÉCELER
LA MALADIE : SI UNE PERSONNE EST INFECTÉE, LE TEST EST POSITIF DANS **99 %**
DES CAS. PAR AILLEURS, LE TEST PRODUIT AUSSI DES **FAUX POSITIFS** : POUR **2 %**
DES PATIENTS SAINS LE TEST EST POSITIF. SI VOTRE TEST EST POSITIF, QUELLE
PROBABILITÉ AVEZ-VOUS D'ÊTRE EFFECTIVEMENT INFECTÉ ?



NOUS AVONS DEUX ÉVÉNEMENTS À GÉRER :

A : LE PATIENT A LA MALADIE

B : LE PATIENT EST TESTÉ POSITIF

L'INFORMATION SUR L'EFFICACITÉ
DU TEST PEUT S'ÉCRIRE :



$$P(A) = 0,001$$

UN PATIENT SUR 1 000 A LA MALADIE.

$$P(B|A) = 0,99$$

LA PROBABILITÉ D'UN TEST POSITIF, SACHANT QUE LE PATIENT EST MALADE, EST DE 0,99.

$$P(B|\text{NON } A) = 0,02$$

LA PROBABILITÉ D'UN FAUX POSITIF, SACHANT QUE LE PATIENT EST SAIN, EST DE 0,02.

ET LA QUESTION EST :

$$P(A|B) = \text{QUOI?}$$

LA PROBABILITÉ D'AVOIR LA MALADIE, SACHANT QUE LE TEST EST POSITIF, EST INCONNUE.

COMME LE TRAITEMENT DE LA MALADIE A DE GRAVES EFFETS SECONDAIRES, LE MÉDECIN, SON AVOCAT ET L'AVOCAT DE SON AVOCAT FONT APPEL À JOE BAYES, CS (CONSULTANT STATISTICIEN) POUR OBTENIR UNE RÉPONSE. JOE UTILISE UN THÉORÈME PROUVÉ PAR SON ANCÊTRE, LE PASTEUR **THOMAS BAYES** (1701-1761).



JOE COMMENCE PAR FAIRE UNE MATRICE 2×2 , QUI DIVISE L'ESPACE ÉCHANTILLON EN 4 CAS MUTUELLEMENT EXCLUSIFS. CHAQUE CAS EST UNE COMBINAISON ENTRE L'ÉTAT DU PATIENT ET LE RÉSULTAT DU TEST.

	A	NON A
B	A ET B	NON A ET B
NON B	A ET NON B	NON A ET NON B

INSCRIVONS LES PROBABILITÉS DE CHAQUE ÉVÉNEMENT DANS LA MATRICE :

	A	NON A	SOMME
B	$P(A \text{ ET } B)$	$P(\text{NON } A \text{ ET } B)$	$P(B)$
NON B	$P(A \text{ ET NON } B)$	$P(\text{NON } A \text{ ET NON } B)$	$P(\text{NON } B)$
SOMME	$P(A)$	$P(\text{NON } A)$	1

LES PROBABILITÉS EN FIN DE LIGNE ET DE COLONNE SONT OBTENUES EN SOMMANT LES LIGNES ET COLONNES.

MAINTENANT, CALCULONS :

$$P(A \text{ ET } B) = P(B|A) P(A) = (0,99)(0,001) = 0,00099$$

$$P(\text{NON } A \text{ ET } B) = P(B|\text{NON } A) P(\text{NON } A) = (0,02)(0,999) = 0,01998$$



CELA NOUS PERMET DE COMPLÉTER LA MATRICE :

	A	NON A	SOMME
B	0,00099	0,01998	0,02097
NON B	$P(A \text{ ET NON } B)$	$P(\text{NON } A \text{ ET NON } B)$	$P(\text{NON } B)$
SOMME	0,001	0,999	1

NOUS POUVONS TROUVER LES PROBABILITÉS MANQUANTES PAR SOUSTRACTION DANS LES COLONNES PUIS, EN ADDITIONNANT DANS LES LIGNES.

LA MATRICE FINALE EST :

	A	NON A	SOMME
B	0,00099	0,01998	0,02097
NON B	0,00001	0,97902	0,97903
SOMME	0,001	0,999	1
	P(A)	P(NON A)	

NOUS POUVONS ALORS EN DÉDUIRE :

$$P(A|B) = \frac{P(A \text{ ET } B)}{P(B)} = \frac{0,00099}{0,02097} = 0,0472$$

MALGRÉ L'APPARENTE FIABILITÉ DU TEST, **MOINS DE 5 %** DE CEUX QUI SONT TESTÉS POSITIFS SONT VÉRITABLEMENT MALADES ! ON APPELLE CELA LE **PARADOXE DES FAUX POSITIFS**.

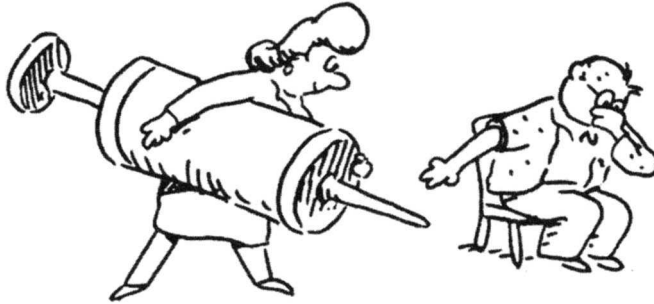
PARADOXE ET PAIRE
D'AVOCATS...



CE TABLEAU EXPLIQUE
LE PROBLÈME AVEC
1000 PATIENTS. EN MOYENNE,
SEULEMENT 21 PERSONNES
SERONT TESTÉES POSITIVES
– ET SEULEMENT **UNE**
PERSONNE SERA EFFECTIVEMENT
MALADE ! IL Y AURA DONC 20
FAUX POSITIFS, CELA PROVIENT
DU FAIT QUE **LE GROUPE**
SAIN EST DE LOIN
LE PLUS IMPORTANT.

	MALADES	BIEN PORTANTS	TOTAL
TESTS POSITIFS	1	20	21
TESTS NÉGATIFS	0	979	979
	1	999	1000

QUE DOIT FAIRE LE MÉDECIN ? JOE BAYES LUI CONSEILLE DE NE PAS COMMENCER LE TRAITEMENT SUR LA SEULE BASE DU TEST. LE TEST RESTE INFORMATIF : AVEC UN RÉSULTAT POSITIF, LA PROBABILITÉ QUE LE PATIENT SOIT MALADE PASSE DE 1 POUR 1000 À 1 POUR 21, MAIS DANS CE CAS LE MÉDECIN DEVRA FAIRE DE **NOUVEAUX TESTS**.



JOE BAYES TOUCHE SON CHÈQUE DE CONSULTANT AVANT D'ADMETTRE QUE TOUTE SON ÉTUDE POUVAIT SE RÉSUMER EN UNE SEULE FORMULE APPELÉE **THÉORÈME DE BAYES**.

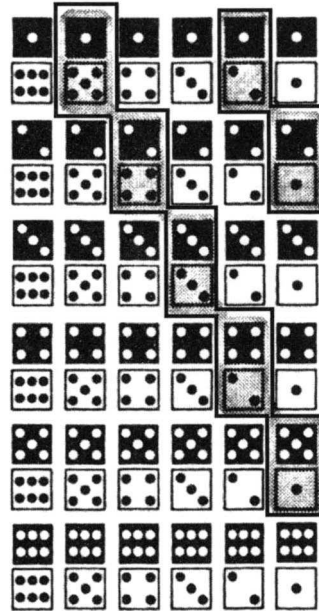
$$P(A|B) = \frac{P(A)P(B|A)}{P(A)P(B|A) + P(\text{NON } A)P(B|\text{NON } A)}$$



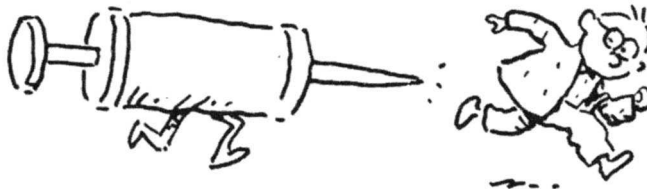
LA FORMULE EXPRIME $P(A|B)$ À PARTIR DE $P(A)$ ET DES DEUX PROBABILITÉS CONDITIONNELLES $P(B|A)$ ET $P(B|\text{NON } A)$. ON LA DÉMONTRE FACILEMENT EN NOTANT QUE LA FRACTION DE DROITE PEUT S'EXPRIMER COMME :

$$\frac{P(A \text{ ET } B)}{P(A \text{ ET } B) + P((\text{NON } A) \text{ ET } B)} = \frac{P(A \text{ ET } B)}{P(B)} = P(A|B)$$

DANS CE CHAPITRE, NOUS AVONS VU LES BASES DES PROBABILITÉS : LES DÉFINITIONS, L'ESPACE ÉCHANTILLON ET LES ÉVÉNEMENTS ÉLÉMENTAIRES, LES PROBABILITÉS CONDITIONNELLES ET CERTAINES FORMULES ESSENTIELLES POUR CALCULER LES PROBABILITÉS. NOUS AVONS ILLUSTRÉ CES IDÉES AVEC L'EXEMPLE D'UN LANCER DE DEUX DÉS. POUR UN JOUEUR MODERNE, LES PROBABILITÉS SONT UN OUTIL DE CHOIX.



ENFIN, AVEC L'EXEMPLE MÉDICAL, NOUS AVONS MONTRÉ COMMENT DES IDÉES ABSTRAITES POUVAIENT NOUS AIDER À PRENDRE DE BONNES DÉCISIONS DANS UN CONTEXTE D'INFORMATION IMPARFAITE ET DE RISQUES RÉELS – CE QUI EST LE BUT ULTIME DES STATISTIQUES.



MAIS CE N'EST QU'UN DÉBUT. POUR NOUS, LES PROBABILITÉS CONSTITUENT UN **OUTIL** (SANS CONTESTE ESSENTIEL) DANS L'ÉTUDE DES STATISTIQUES. DANS LES CHAPITRES QUI SUIVENT, NOUS ALLONS EXPLORER LA RELATION SUBTILE ENTRE LES PROBABILITÉS, LES VARIATIONS DANS LES DONNÉES STATISTIQUES, ET NOTRE CONFIANCE DANS L'INTERPRÉTATION DE NOS OBSERVATIONS.





Chapitre 4

LES VARIABLES ALÉATOIRES

DANS LE CHAPITRE 2, NOUS AVONS VU QUE L'OBSERVATION DE DONNÉES NUMÉRIQUES, COMME LE POIDS DES ÉTUDIANTS, PEUT ÊTRE REPRÉSENTÉE GRAPHIQUEMENT OU RÉSUMÉE EN TERMES DE TENDANCE CENTRALE, DE DISPERSION, DE VALEURS EXTRÊMES, ETC.

DANS LE CHAPITRE 3, NOUS AVONS VU COMMENT DES PROBABILITÉS PEUVENT ÊTRE AFFECTÉES À DES RÉSULTATS D'UNE EXPÉRIENCE ALÉATOIRE.



SI NOUS IMAGINONS UNE EXPÉRIENCE ALÉATOIRE RÉPÉTÉE DE NOMBREUSES FOIS, NOUS NOUS ATTENDONS À CE QUE LES FRÉQUENCES DES RÉSULTATS OBSERVÉS TENDENT VERS LEURS PROBABILITÉS À TERME. LES PROBABILITÉS CONSTITUENT UN **MODÈLE** POUR LES EXPÉRIENCES DE LA VIE RÉELLE... ALORS, POURQUOI NE PAS FAIRE L'ÉTUDE D'UN MODÈLE À PARTIR DE L'ANALYSE DES DONNÉES DÉCRITES PAR CE MODÈLE ?

LE CONCEPT ESSENTIEL EST CELUI DE VARIABLE ALÉATOIRE, QUE NOUS ÉCRIVONS AVEC UN GRAND



X

UNE VARIABLE ALÉATOIRE EST DÉFINIE COMME LE **RÉSULTAT NUMÉRIQUE** D'UNE **EXPÉRIENCE ALÉATOIRE**.

PAR EXEMPLE, IMAGINONS QU'ON SÉLECTIONNE UN ÉTUDIANT AU HASARD DE NOTRE GROUPE D'ÉTUDIANTS. IL S'AGIT BIEN D'UNE EXPÉRIENCE ALÉATOIRE. LA **TAILLE**, LE **POIDS**, LE **REVENU FAMILIAL**, LE **RÉSULTAT AU BAC** ET LA **MOYENNE GÉNÉRALE** DE L'ÉTUDIANT SONT DES **VARIABLES NUMÉRIQUES** DÉCRIVANT LES PROPRIÉTÉS DE L'ÉTUDIANT TIRÉ AU HASARD. CE SONT TOUTES DES VARIABLES ALÉATOIRES.



UN AUTRE EXEMPLE : LANCER DEUX PIÈCES (L'EXPÉRIENCE ALÉATOIRE) ET NOTER LE **NOMBRE** DE FACES : 0, 1 OU 2.

RÉSULTAT

x

PP

|

0

PF

FP

\

1

FF

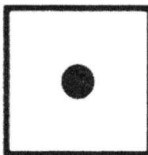
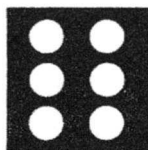
|

2



ATTENTION POUR LA NOTATION, LA VARIABLE S'ÉCRIT AVEC UN **X** MAJUSCULE. LE x EN MINUSCULE REPRÉSENTE UNE VALEUR PARTICULIÈRE DE X , AINSI $x = 2$ VEUT DIRE QU'IL Y A EU DEUX FOIS FACE.

VOICI UN AUTRE EXEMPLE BASÉ
SUR LE LANCER CLASSIQUE DE DEUX
DÉS. SOIT Y LA SOMME DES POINTS
VISIBLES SUR LES DEUX DÉS.
POUR CETTE VARIABLE ALÉATOIRE,
 Y EST UN NOMBRE COMPRIS
ENTRE 2 ET 12.



$$Y = 7$$

NOUS VOULONS MAINTENANT EXAMINER LES **PROBABILITÉS** DES RÉSULTATS.
ON ÉCRIT $P(X = x)$, OU SIMPLEMENT $P(x)$, LA PROBABILITÉ QUE LA VARIABLE
ALÉATOIRE X SOIT ÉGALE À x . POUR LE CAS DE LA VARIABLE ALÉATOIRE DU LANCER
DE DEUX PIÈCES, NOUS AVONS LE TABLEAU SUIVANT :

x	0	1	2
$P(X = x)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

CE TABLEAU EST APPELÉ
LA **DISTRIBUTION**
DE PROBABILITÉS
DE LA VARIABLE ALÉATOIRE X .

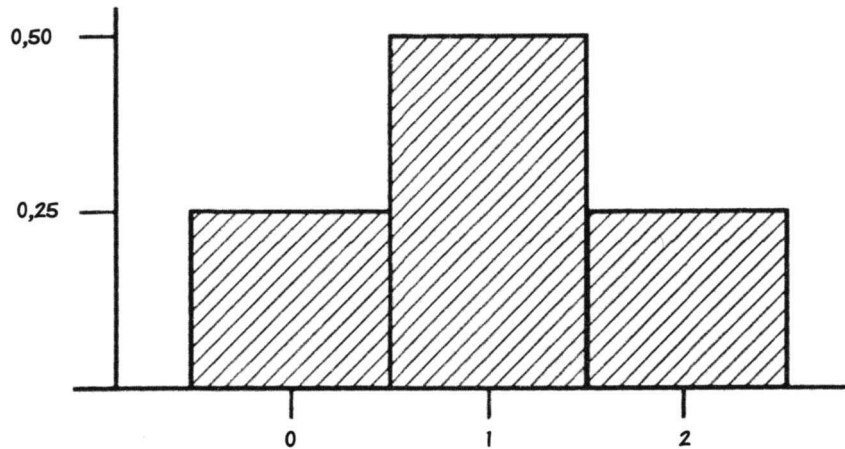
POUR LA VARIABLE ALÉATOIRE Y (SOMME DU LANCER DE DEUX DÉS), LA DISTRIBUTION
DE PROBABILITÉS RESSEMBLE À CECI :

y	2	3	4	5	6	7	8	9	10	11	12
$P(Y = y)$	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$



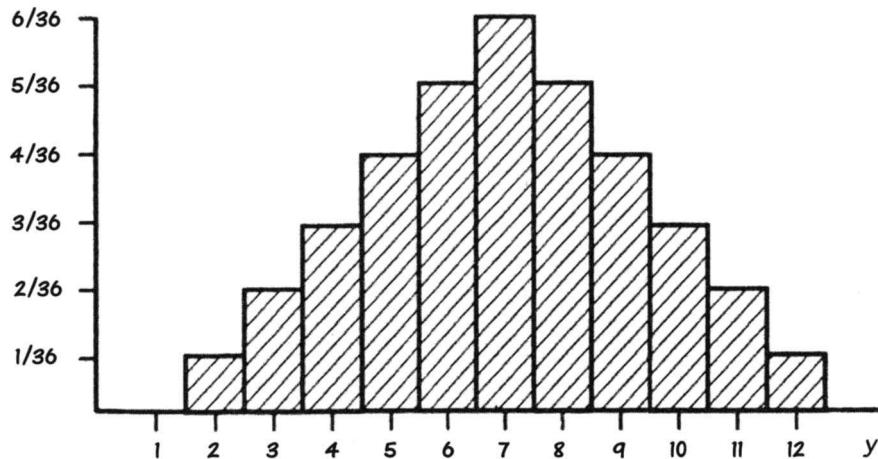
OUEP! C'EST
POUR CELA QUE
J'AI ARRÊTÉ
LES DÉS!

MAINTENANT DESSINONS DES GRAPHIQUES, OU DES **HISTOGRAMMES** REPRÉSENTANT LES DISTRIBUTIONS DE PROBABILITÉS. POUR CHAQUE VALEUR DE X , ON DESSINE UNE BARRE DE HAUTEUR $p(x)$.

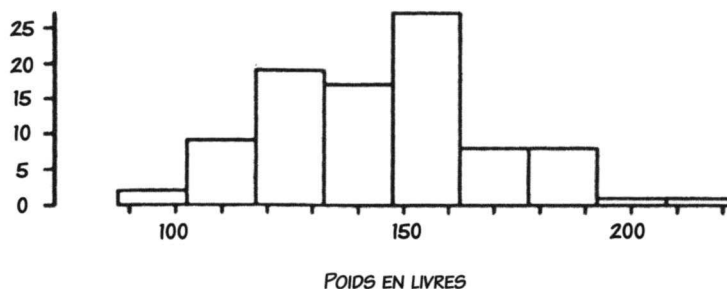


IL EST FACILE DE VOIR QUE L'AIRE TOTALE DE CES RECTANGLES EST 1. CHAQUE RECTANGLE A POUR BASE 1 ET UNE HAUTEUR DE $p(x)$. AINSI L'AIRE TOTALE EST ÉGALE À LA SOMME DES PROBABILITÉS DE TOUS LES RÉSULTATS POSSIBLES, C'EST-À-DIRE 1.

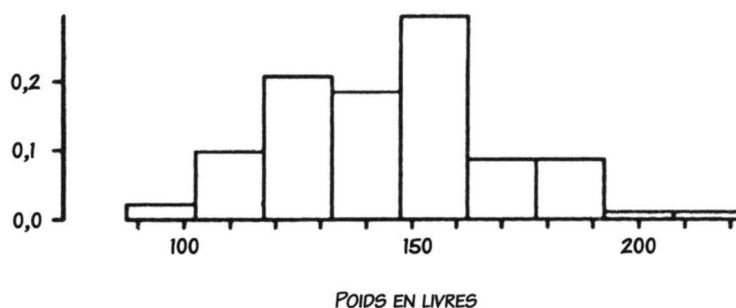
VOICI L'HISTOGRAMME DES PROBABILITÉS DE LA VARIABLE ALÉATOIRE Y :
IL DÉCRIT LA DISTRIBUTION DE PROBABILITÉS DE LA SOMME DE DEUX DÉS.



POURQUOI APPELLE-T-ON CES GRAPHIQUES DES HISTOGRAMMES ? VOUS VOUS RAPPELEZ QUE DANS LE CHAPITRE 2, UN HISTOGRAMME ÉTAIT UN GRAPHIQUE QUI REPRÉSENTAIT LE NOMBRE DE DONNÉES OBSERVÉES DANS CHACUNE DES CLASSES.



À PARTIR DE CET HISTOGRAMME DES **EFFECTIFS**, ON A DÉFINI UN HISTOGRAMME DES **FRÉQUENCES** QUI REPRÉSENTE LA **PROPORTION** DE DONNÉES OBSERVÉES DANS CHAQUE CLASSE.



MAIS, VOUS VOUS SOUVENEZ QUE PAR DÉFINITION UNE PROBABILITÉ REPRÉSENTE LA FRÉQUENCE D'UN ÉVÉNEMENT SUR LE « LONG TERME ». SI L'ON RÉPÈTE L'EXPÉRIENCE ALÉATOIRE DE NOMBREUSES FOIS, L'HISTOGRAMME EN **FRÉQUENCE** DES RÉSULTATS OBSERVÉS DOIT ÊTRE SEMBLABLE À L'HISTOGRAMME DES **PROBABILITÉS** DE LA VARIABLE ALÉATOIRE.



ILLUSTRONS CELA AVEC LA VARIABLE ALÉATOIRE X ET UNE FURIEUSE LANCEUSE DE PIÈCES.



LA LANCEUSE COMMENCE À LANCER DEUX PIÈCES DE FAÇON RÉPÉTITIVE, EN CONSERVANT LE RÉSULTAT CHAQUE FOIS.



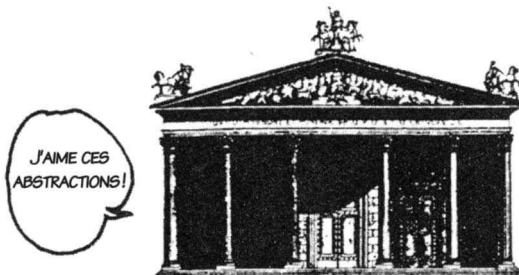
NOUS CONNAISSONS LA DISTRIBUTION DE PROBABILITÉS DE X ET NOUS SAVONS QUE LES RÉSULTATS DE L'EXPÉRIENCE RÉELLE CORRESPONDENT APPROXIMATIVEMENT AUX PROBABILITÉS. APRÈS 1000 LANCERS, LA LANCEUSE FAIT LE POINT SUR SES DONNÉES :

MODÈLE DE PROBABILITÉS		DONNÉES OBSERVÉES	
$p(x)$	x	n_x = NOMBRE D'OCCURRENCES	n_x/n = FRÉQUENCE
0,25	0	260	0,260
0,5	1	517	0,517
0,25	2	223	0,223

ET NOUS CONSTATONS QUE L'HISTOGRAMME DES PROBABILITÉS DE X CORRESPOND À UNE « FORME PURE » OU AU MODÈLE DE L'HISTOGRAMME EN FRÉQUENCE DES DONNÉES OBSERVÉES.



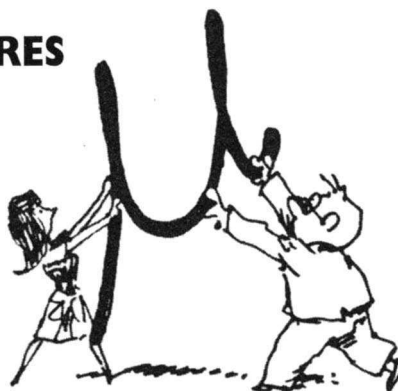
POUR POURSUIVRE L'ANALOGIE ENTRE FRÉQUENCE ET DONNÉES, NOUS ALLONS PARLER DE LA MOYENNE ET DE LA VARIANCE (OU DE L'ÉCART-TYPE) D'UNE DISTRIBUTION DE PROBABILITÉS...



ET AFIN DE NOUS RAPPELER QUE NOUS SOMMES DANS LE DOMAINE DE L'ABSTRACTION, NOUS MOBILISONS DES LETTRES GRECQUES...

MOYENNE et VARIANCE des VARIABLES ALÉATOIRES

ON UTILISE UNE TERMINOLOGIE ET DES SYMBOLES PARTICULIERS POUR DISTINGUER LES CARACTÉRISTIQUES DE DONNÉES OBSERVÉES DE CELLES DE PROBABILITÉS DE DISTRIBUTION.



LES CARACTÉRISTIQUES NUMÉRIQUES DE DONNÉES SONT APPELÉES CARACTÉRISTIQUES D'**ÉCHANTILLON**, ALORS QUE LES CARACTÉRISTIQUES D'UNE PROBABILITÉ DE DISTRIBUTION SONT APPELÉES CARACTÉRISTIQUES DU **MODÈLE** OU DE LA **POPULATION**. ON UTILISE DES LETTRES GRECQUES COMME μ (MU) POUR LA MOYENNE DE POPULATION ET σ (SIGMA MINUSCULE) POUR L'ÉCART-TYPE DE POPULATION (POUR UN ÉCHANTILLON DE DONNÉES, ON UTILISE LES LETTRES ROMAINES \bar{x} ET s).



PARCE QUE
LES ROMAINS ÉTAIENT
LÉGERS POUR LA THÉORIE
MAIS FORTS
EN CONSTRUCTION...
ET D'AUTRES CHOSES
DE CE GENRE.

LA MOYENNE ÉCHANTILLON
A ÉTÉ DÉFINIE PAR L'ÉQUATION

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$



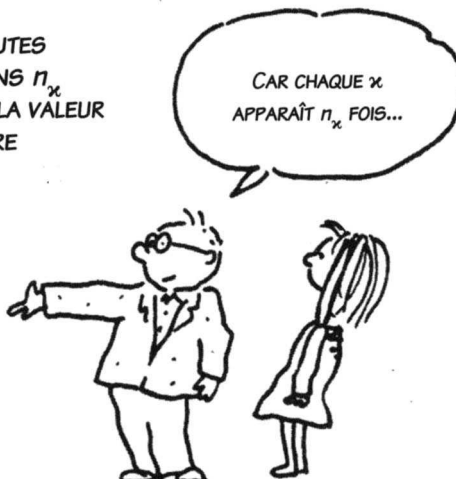
CERTAINES DE CES DONNÉES OBSERVÉES x_i PEUVENT AVOIR DES VALEURS ÉGALES. DANS L'EXEMPLE DE LA LANCEUSE DE PIÈCES, LES SEULES VALEURS POSSIBLES SONT 0, 1 OU 2 ET CELLE-CI A FAIT 1000 LANCERS. LA VALEUR 0 EST APPARUE 260 FOIS, LA VALEUR 1 : 517 FOIS, ET LA VALEUR 2 : 223 FOIS.

LA VARIABLE x VA PRENDRE TOUTES
LES VALEURS POSSIBLES, NOTONS n_x
LE NOMBRE DE DONNÉES DONT LA VALEUR
EST x . ON PEUT ALORS RÉÉCRIRE
LA FORMULE :

$$\bar{x} = \frac{1}{n} \sum_{\text{TOUT } x} n_x x$$

OU

$$\bar{x} = \sum_{\text{TOUT } x} x \frac{n_x}{n}$$



AH! MAIS n_x/n CORRESPOND À LA **FRÉQUENCE**... C'EST-À-DIRE LA «PROBABILITÉ APPROXIMÉE...» SOIT LE NOMBRE QUI TEND VERS $p(x)$... DONC PAR ANALOGIE NOUS OBTENONS LA FORMULE

$$\sum_{\text{TOUT } x} x p(x)$$



ET DÉFINISSONS CELA COMME
L'**ESPÉRANCE** (OU MOYENNE)
DE LA **DISTRIBUTION DE PROBABILITÉ**.

Définition : LA moyenne d'une VARIABLE ALÉATOIRE X EST DÉFINIE PAR :

$$\mu = \sum_{\text{TOUT } x} xp(x)$$

ET C'EST
LE CENTRE DE CET
HISTOGRAMME !



ON L'APPELLE AUSSI ESPÉRANCE OU **VALEUR ESPÉRÉE** DE X , OU $E[X]$.
ON PEUT RETENIR QU'IL S'AGIT DE LA SOMME DE TOUTES LES VALEURS POSSIBLES
PONDÉRÉE PAR LES PROBABILITÉS.

DANS L'EXPÉRIENCE DE LA LANCEUSE DE PIÈCES, ON PEUT COMPARER LA MOYENNE
D'ÉCHANTILLON \bar{x} AVEC LA MOYENNE DE POPULATION μ :

ÉCHANTILLON		
x	$\frac{n_x}{n}$	$x \frac{n_x}{n}$
0	0,26	0,0
1	0,517	0,517
2	0,223	0,446
		<u>0,963 = \bar{x}</u>

MODÈLE		
x	$p(x)$	$xp(x)$
0	0,25	0,0
1	0,5	0,5
2	0,25	0,5
		<u>1 = μ</u>

MAINTENANT, FAISONS LA MÊME CHOSE
AVEC LA **VARIANCE**. PEUT-ÊTRE VOUS
RAPPELEZ-VOUS LA FORMULE :

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

CELA MESURE (PRESQUE) LA MOYENNE
DES CARRÉS DES ÉCARTS À LA MOYENNE.
COMME PRÉCÉDEMMENT, ON PEUT
LA RÉÉCRIRE :

$$s^2 = \sum_{\text{TOUT } x} (x_i - \bar{x})^2 \frac{n_x}{n-1}$$



IL S'AGIT BIEN D'UNE SOMME PONDÉRÉE D'ÉCARTS AU CARRÉ, MIS À PART CE $n - 1$ AU DÉNOMINATEUR AU LIEU DE n ... NOUS DÉFINISSONS DONC :

LA **variance** D'UNE VARIABLE ALÉATOIRE X COMME LA VALEUR ESPÉRÉE DES CARRÉS DES ÉCARTS À LA MOYENNE DE POPULATION :

$$\sigma^2 = \sum_{\text{TOUT } x} (x - \mu)^2 p(x)$$

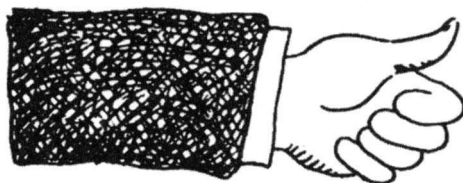
L'**écart-type** σ EST LA RACINE CARRÉE DE LA VARIANCE.

AVEZ-VOUS REMARQUÉ
QUE σ^2 EST AUSSI
 $E[(x - \mu)^2]$?



POUR TROUVER LA VARIANCE
DU CAS DU LANCER DE 2 PIÈCES
(POUR LEQUEL $\mu = 1$),
ON UTILISE LE TABLEAU
DE LA PAGE PRÉCÉDENTE.

x	$p(x)$	$(x - \mu)^2 p(x)$
0	0,25	$(0 - 1)^2 (0,25) = 0,25$
1	0,5	$(1 - 1)^2 (0,5) = 0$
2	0,25	$(2 - 1)^2 (0,25) = 0,25$
		$0,50 = \sigma^2$



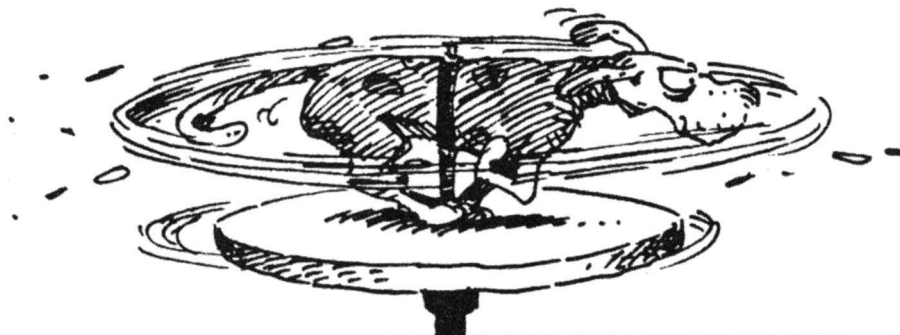
EN RÉSUMÉ : μ ET σ , LES MOYENNES ET ÉCARTS-TYPES DE POPULATION SONT DES CARACTÉRISTIQUES QUE L'ON CALCULE POUR DES **DISTRIBUTIONS DE PROBABILITÉS**. ELLES SONT ANALOGUES AUX MOYENNES \bar{x} ET AUX ÉCARTS-TYPES s D'ÉCHANTILLON QUE L'ON CALCULE À PARTIR DE DONNÉES D'ÉCHANTILLON.

JUSQU'À PRÉSENT NOS EXEMPLES DE VARIABLES ALÉATOIRES ÉTAIENT DISCRETS. LES RÉSULTATS ÉTAIENT DES VALEURS ISOLÉES (OU « DISCRÈTES ») COMME DANS LE CHAPITRE 3, MAIS IL Y A AUSSI DES **variables aléatoires continues**.

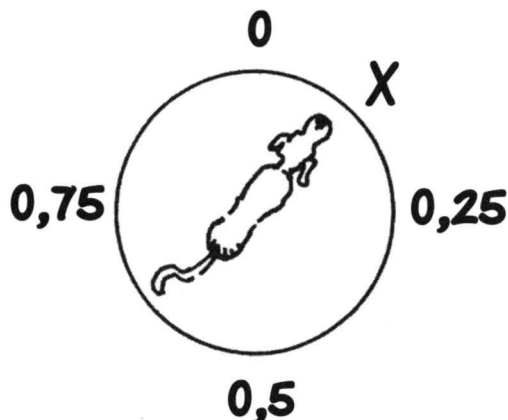
IMAGINONS UNE EXPÉRIENCE ALÉATOIRE OÙ CHAQUE RÉSULTAT AIT UNE PROBABILITÉ DE ZÉRO. AUTREMENT DIT, $p(x) = 0$ POUR TOUT x .



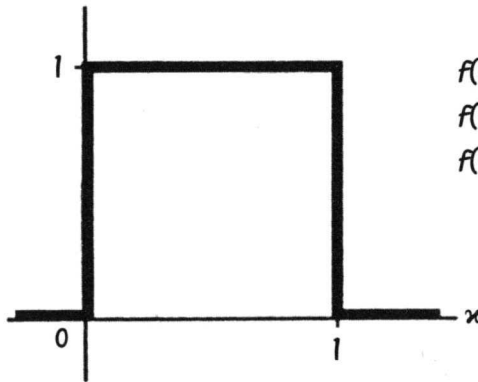
UN EXEMPLE SIMPLE EST CELUI D'UN POINTEUR SUR UNE ROUE (SI VOUS N'AVEZ PAS DE CHIEN DE CHASSE, UNE AIGUILLE FERA L'AFFAIRE). CE DERNIER PEUT S'ARRÊTER N'IMPORTE OÙ DANS LE CERCLE. SOIT X LA PROPORTION DE LA CIRCONFÉRENCE ATTEINTE PAR LE POINTEUR. LA VARIABLE ALÉATOIRE X PEUT PRENDRE N'IMPORTE QUELLE VALEUR ENTRE 0 ET 1, ET DONC UN ÉVENTAIL INFINI DE VALEURS.



IL EST FACILE DE TROUVER LA PROBABILITÉ QUE X APPARTIENNE À UN INTERVALLE DONNÉ :
 PAR EXEMPLE, $P(0,25 \leq X \leq 0,75) = 0,5$
 CAR IL S'AGIT DE LA MOITIÉ DU CERCLE.
 MAIS QU'EN EST-IL DE $P(X = 0,5)$?
 COMME X PEUT PRENDRE UN NOMBRE INFINI DE VALEURS ET QUE CHAQUE VALEUR EST ÉQUIPROBABLE, LA PROBABILITÉ QUE X SOIT **EXACTEMENT 0,5** (OU TOUTE AUTRE VALEUR) EST PRÉCISÉMENT NULLE.



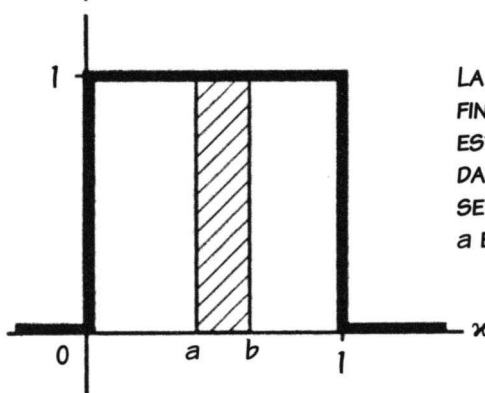
COMMENT REPRÉSENTER CELA ?
 PAR ANALOGIE AVEC LE CAS DES PROBABILITÉS
 DISCRÈTES, IL FAUT VOIR LES PROBABILITÉS
 CONTINUES COMME L'AIRE SITUÉE SOUS
QUELQUE CHOSE. DANS LE CAS DU POINTEUR
 DE LA ROUE, CE QUELQUE CHOSE RESSEMBLE
 À CECI :



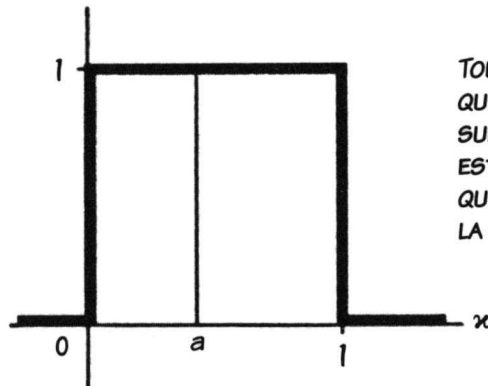
$$f(x) = 0 \text{ si } x < 0$$

$$f(x) = 1 \text{ si } 0 \leq x \leq 1$$

$$f(x) = 0 \text{ si } x > 1$$



LA PROBABILITÉ QUE LE POINTEUR
 FINISSE SA COURSE ENTRE a ET b
 EST PRÉCISÉMENT L'AIRE SITUÉE
 DANS LA PARTIE HACHURÉE QUI
 SE TROUVE SOUS LA COURBE ENTRE
 a ET b . CETTE AIRE VAUT $b - a$.



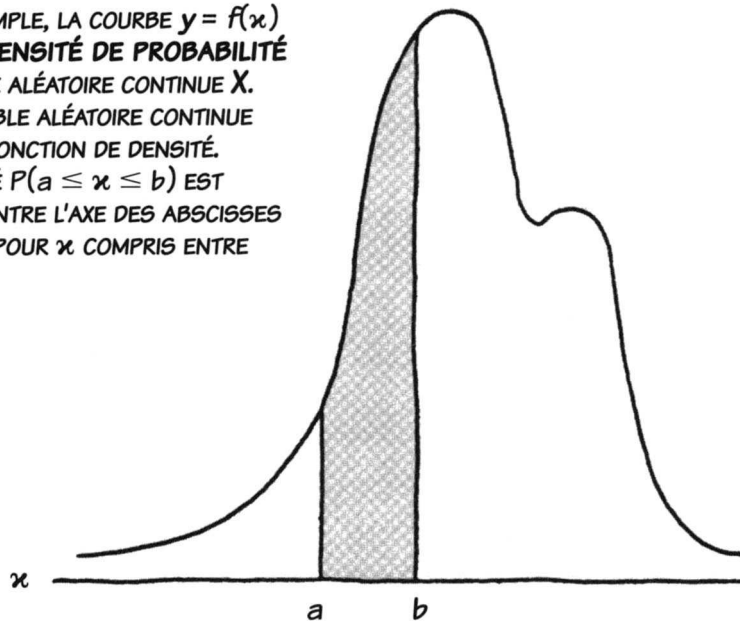
TOUTEFOIS, LA PROBABILITÉ
 QUE LE POINTEUR TOMBE
 SUR UNE VALEUR FIXE
 EST DE **ZÉRO** (NOTEZ AUSSI
 QUE L'AIRE TOTALE SOUS
 LA COURBE VAUT EXACTEMENT 1).

LES NOMBRES ALÉATOIRES GÉNÉRÉS PAR ORDINATEUR OU AVEC DES CALCULATRICES ONT LE MÊME TYPE DE REPRÉSENTATION. IL SUFFIT DE PRESSER UN BOUTON, ET HOP! ON OBTIENT UN NOMBRE ENTRE 0 ET 1. CHAQUE NOMBRE EST ÉQUIPROBABLE COMME DANS LE CAS DU POINTEUR DE LA ROUE.



MAIS MALHEUREUSEMENT, CES NOMBRES NE SONT PAS VRAIMENT ALÉATOIRES. ILS SONT PRODUITS PAR DES ALGORITHMES. ON PARLE PLUS PRÉCISÉMENT DE NOMBRES PSEUDO-ALÉATOIRES.

DANS CET EXEMPLE, LA COURBE $y = f(x)$ EST APPELÉE **DENSITÉ DE PROBABILITÉ** DE LA VARIABLE ALÉATOIRE CONTINUE X . CHAQUE VARIABLE ALÉATOIRE CONTINUE A SA PROPRE FONCTION DE DENSITÉ. LA PROBABILITÉ $P(a \leq x \leq b)$ EST ALORS L'AIRE ENTRE L'AXE DES ABSCISSES ET LA COURBE POUR x COMPRIS ENTRE a ET b .



EN GÉNÉRAL, LA FONCTION DE DENSITÉ N'EST PAS SI SIMPLE, ET LE CALCUL DE L'AIRE EST LOIN D'ÊTRE TRIVIAL.



ON UTILISE UNE NOTATION MATHÉMATIQUE POUR DÉCRIRE CETTE AIRE SOUS LA FONCTION $f(x)$. CE SYMBOLE SE LIT «L'INTÉGRALE DE f ENTRE a ET b ».

$$\int_a^b f(x) dx$$



COMME POUR LES PROBABILITÉS DISCRÈTES, LES FONCTIONS DE DENSITÉS CONTINUES VÉRIFIENT DEUX PROPRIÉTÉS :

$$f(x) \geq 0$$

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

(NE VOUS AFFOLEZ PAS DE LA PRÉSENCE DE L'INFINI. ON VEUT JUSTE DIRE QUE L'ON CHERCHE L'AIRE TOTALE SOUS LA COURBE D'UN BOUT À L'AUTRE.)



BIEN QUE LA NOTATION
PUISSE SEMBLER
ÉTRANGE, ELLE SIGNIFIE
SIMPLEMENT UNE AIRE...
LE SIGNE D'INTÉGRATION
EST LUI-MÊME UN S
ALLONGÉ SIGNIFIANT
S COMME SOMME.



EN UTILISANT L'INTÉGRALE À LA PLACE DE LA SOMME, ON DÉFINIT LA **MOYENNE**
et la **VARIANCE d'une variable aléatoire continue.**

$$\mu = \int_{-\infty}^{\infty} xf(x)dx$$

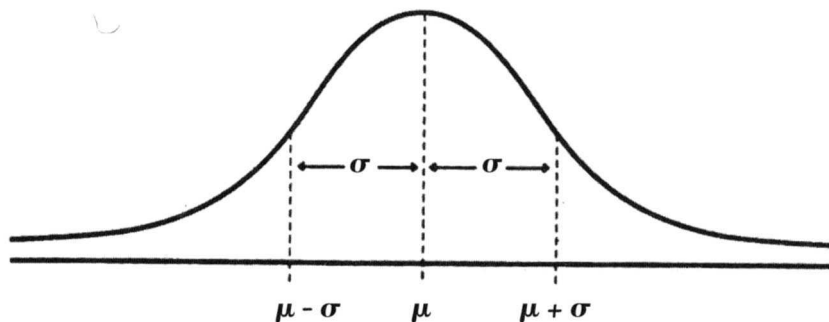
PAR ANALOGIE
AVEC LES
FORMULES
DISCRÈTES :

$$\mu = \sum_{\text{TOUT } x} xp(x)$$

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x)dx$$

$$\sigma^2 = \sum_{\text{TOUT } x} (x - \mu)^2 p(x)$$

BIEN QUE CELA NE SE VOIE PAS DIRECTEMENT DANS LES FORMULES,
LES DÉFINITIONS DE MOYENNE ET DE VARIANCE SONT TOTALEMENT COHÉRENTES
AVEC L'INTERPRÉTATION DE TENDANCE CENTRALE ET DE DISPERSION MOYENNE
DES PROBABILITÉS DÉFINIES PAR LA DENSITÉ $f(x)$. LE GRAPHIQUE À AVOIR EN TÊTE
EST LE SUIVANT :



SOMMATION de variables aléatoires

UNE FOIS CONNUES LA MOYENNE
ET LA VARIANCE D'UNE VARIABLE ALÉATOIRE,
QUE PEUT-ON EN FAIRE? TOUT D'ABORD,
ON PEUT EN DÉDUIRE LA MOYENNE
ET LA VARIANCE D'AUTRES VARIABLES ALÉATOIRES.



PAR EXEMPLE, REGARDONS LE LANCER D'UNE PIÈCE. POSONS $X = 1$ SI LA PIÈCE
EST FACE ET 0 SI ELLE EST PILE.

x	0	1
$p(x)$	0,5	0,5

À PRÉSENT, VOUS DEVRIEZ POUVOIR
CALCULER LA MOYENNE :

$$\begin{aligned} E[x] &= 0 \times p(0) + 1 \times p(1) \\ &= 0 + 0,5 \\ &= 0,5 \end{aligned}$$

ET LA VARIANCE :

$$\begin{aligned} \sigma^2 &= (0 - 0,5)^2 p(0) + (1 - 0,5)^2 p(1) \\ &= 0,25 \end{aligned}$$



MAINTENANT JOUONS À UN JEU SIMPLE. D'ABORD VOUS MISEZ 6 € POUR JOUER.
JE LANCE UNE PIÈCE, VOUS GAGNEZ 10 € SI C'EST FACE, ET RIEN SI C'EST PILE.
AINSI, VOS GAINS G SONT :

$$G = 10X - 6$$

C'EST UNE NOUVELLE VARIABLE
ALÉATOIRE! QUELLES SONT
SA MOYENNE ET SA VARIANCE?



UN PEU DE RÉFLEXION DEVRAIT VOUS
CONVAINCRE QUE $E[G]$ VÉRIFIE :

$$\begin{aligned} E[G] &= E[10X - 6] \\ &= 10E[X] - 6 \end{aligned}$$

CE QUI REVIENT À :
 $10(0,5) - 6 = -1$

CE QUE L'ON PEUT VÉRIFIER
AVEC LE TABLEAU :

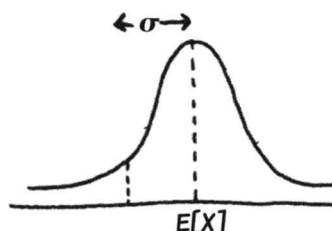
X	0	1
G	-6	4
$p(g)$	0,5	0,5

DONC, VOTRE
« GAIN » ESPÉRÉ
EST UNE PERTE !



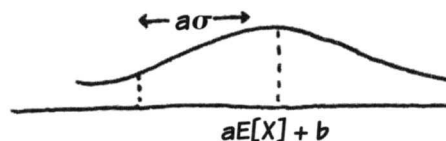
DE FAÇON GÉNÉRALE, IL EST FACILE
DE MONTRER QUE :

$$E[aX + b] = aE[X] + b$$



POUR TOUT NOMBRE a ET b ET X ,
UNE VARIABLE ALÉATOIRE.
POUR LA VARIANCE NOUS AVONS AUSSI
LE RÉSULTAT GÉNÉRAL SUIVANT :

$$\sigma^2(aX + b) = a^2 \sigma^2(X)$$



DANS LE JEU PRÉCÉDENT, LES RÉSULTATS
POSSIBLES SONT -6 ET 4, IL EST DONC
ÉVIDENT QUE LA VARIANCE DE G EST PLUS
IMPORTANTE QUE CELLE DE X . EN EFFET

$$\begin{aligned} \sigma^2(G) &= \sigma^2(10X + 6) \\ &= 100\sigma^2(X) \\ &= 25 \end{aligned}$$

ET

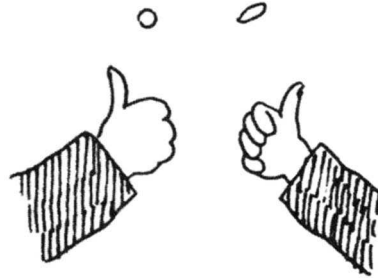
$$\sigma(G) = 5$$



ÇA M'A L'AIR
D'ÊTRE
UN PARI
DE DUPE !

ON PEUT AUSSI AJOUTER DEUX VARIABLES ALÉATOIRES. PAR EXEMPLE, SUPPOSONS QUE L'ON LANCE **DEUX FOIS** UNE PIÈCE. LE NOMBRE DE FACES SUR LES DEUX LANCERS EST $X_1 + X_2$, OÙ X_1 ET X_2 SONT LES VARIABLES ALÉATOIRES DU PREMIER ET DU SECOND LANCER.

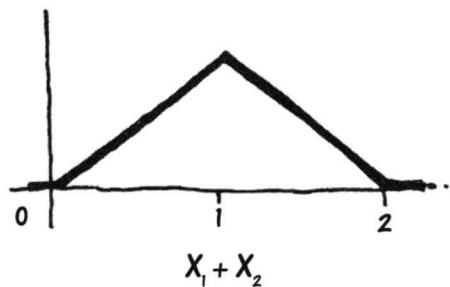
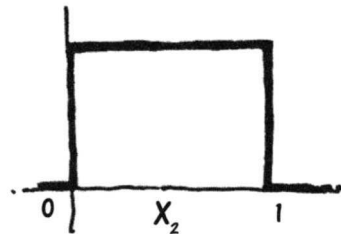
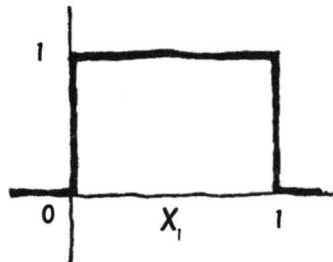
$x_1 + x_2$	0	1	2
$p(x_1 + x_2)$	0,25	0,5	0,25



À NOUVEAU, IL EST FACILE DE VOIR QUE :

$$E[X_1 + X_2] = E[X_1] + E[X_2]$$

(NE DEMANDEZ PAS QUELLE EST LA DISTRIBUTION DE PROBABILITÉ DE $X_1 + X_2$, CAR LE LIEN AVEC LES DISTRIBUTIONS ORIGINALES DE X_1 ET X_2 EST COMPLIQUÉ. PAR EXEMPLE, SI X_1 ET X_2 SONT DEUX LANCERS DE POINTEURS SUR UNE ROUE, LA DENSITÉ DE PROBABILITÉ RESSEMBLE À CECI :)



LA VARIANCE DE LA SOMME DE DEUX VARIABLES ALÉATOIRES A UNE FORME SIMPLE LORSQUE LES VARIABLES X ET Y SONT INDÉPENDANTES. LA DÉFINITION TECHNIQUE DE L'INDÉPENDANCE EST BASÉE SUR LA PROPRIÉTÉ QUE $P(A \text{ ET } B) = P(A)P(B)$... MAIS, POUR NOUS, L'INDÉPENDANCE SIGNIFIE QUE X ET Y SONT GÉNÉRÉS PAR DES PROCESSUS INDÉPENDANTS, COMME LE LANCER DE PIÈCES OU DE DÉS, ETC.



EN DEHORS DU CASINO,
IL EST DIFFICILE DE TROUVER
UNE INDÉPENDANCE
PARFAITE...

LORSQUE X ET Y SONT
INDÉPENDANTS, LEURS
VARIANCES S'ADDITIONNENT :

$$\sigma^2(X + Y) = \sigma^2(X) + \sigma^2(Y)$$

DANS LE CAS DU LANCER DE DEUX PIÈCES :

$$\begin{aligned}\sigma^2(X_1 + X_2) &= \sigma^2(X_1) + \sigma^2(X_2) \\ &= 0,25 + 0,25 \\ &= 0,50\end{aligned}$$



MAIS DANS
LE MONDE IDÉAL
DES STATISTIQUES,
C'EST TRÈS UTILE...

TOUT CELA PEUT SE GÉNÉRALISER À UNE SOMME DE PLUSIEURS VARIABLES ALÉATOIRES :

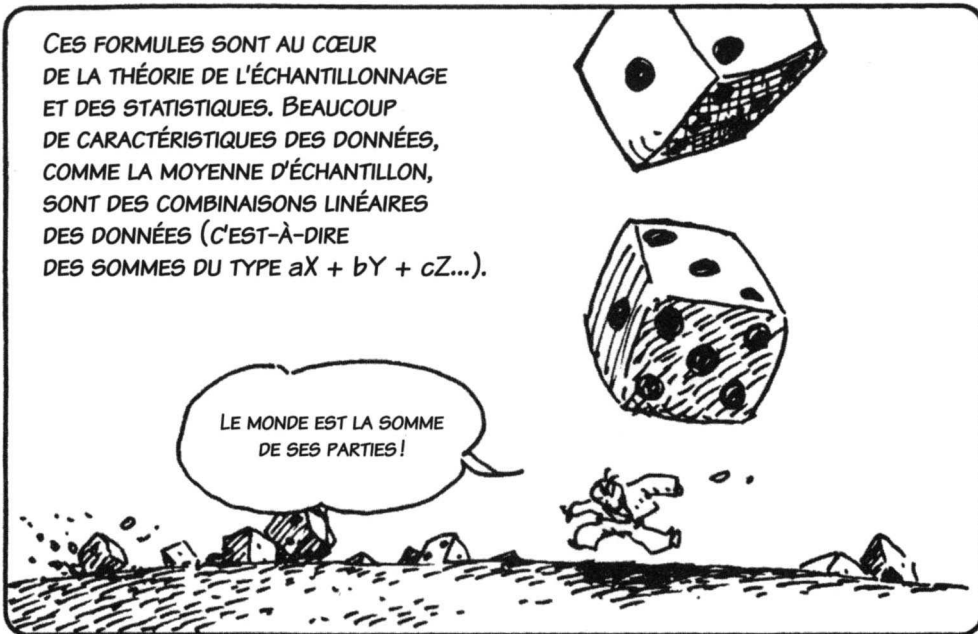
$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i]$$

ET LORSQUE LES X_i SONT DES VARIABLES
ALÉATOIRES INDÉPENDANTES,

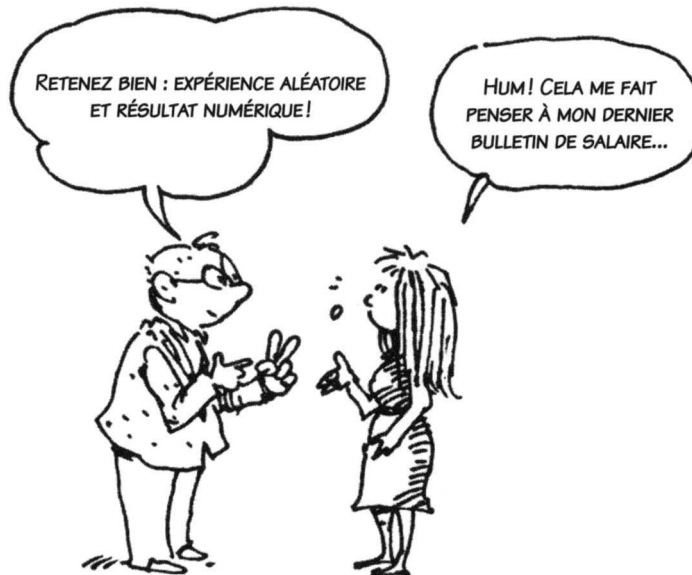
$$\sigma^2\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \sigma^2(X_i)$$



CES FORMULES SONT AU CŒUR DE LA THÉORIE DE L'ÉCHANTILLONNAGE ET DES STATISTIQUES. BEAUCOUP DE CARACTÉRISTIQUES DES DONNÉES, COMME LA MOYENNE D'ÉCHANTILLON, SONT DES COMBINAISONS LINÉAIRES DES DONNÉES (C'EST-À-DIRE DES SOMMES DU TYPE $aX + bY + cZ...$).



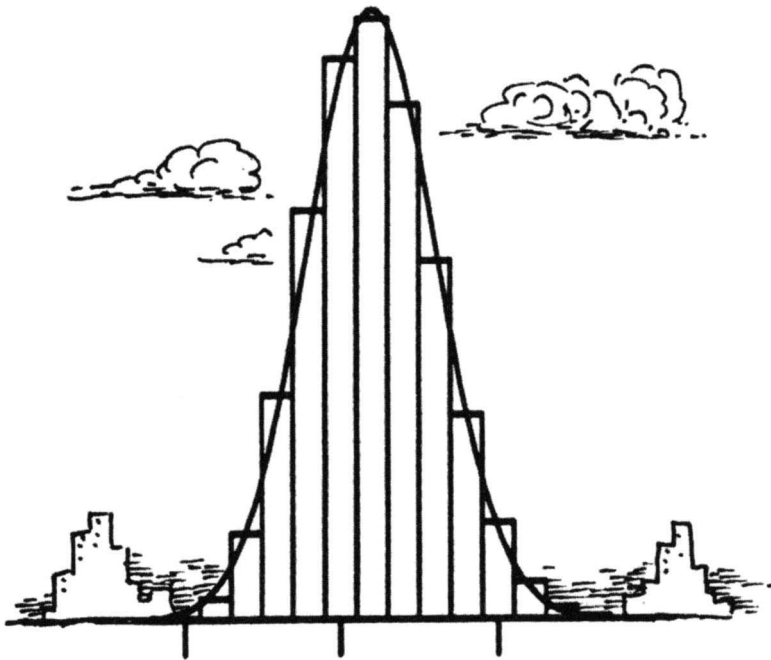
DANS LE PROCHAIN CHAPITRE, NOUS ALLONS VOIR DEUX EXEMPLES IMPORTANTS DE VARIABLES ALÉATOIRES : LE PREMIER, LA **LOI BINOMIALE**, EST UNE SOMME RÉPÉTÉE DE VARIABLES ALÉATOIRES INDÉPENDANTES. LE SECOND, LA **LOI NORMALE**, EST UNE VARIABLE ALÉATOIRE CONTINUE QUI A UN LIEN SURPRENANT AVEC LA LOI BINOMIALE ET AUSSI AVEC D'AUTRES SOMMES DE VARIABLES INDÉPENDANTES.



Chapitre 5

Une histoire de deux distributions

NOUS ALLONS MAINTENANT EXAMINER DEUX EXEMPLES DE VARIABLES ALÉATOIRES, DONT L'UNE EST DISCRÈTE ET L'AUTRE CONTINUE.



NOUS ALLONS COMMENCER PAR LA DISTRIBUTION DISCRÈTE APPELÉE VARIABLE ALÉATOIRE **BINOMIALE** (OU LOI BINOMIALE). SUPPOSONS QUE NOUS AYONS UN PROCESSUS ALÉATOIRE AVEC SEULEMENT **DEUX** RÉSULTATS : UN LANCER DE PIÈCE, UN MATCH SPORTIF, UN CONTRÔLE DE POLLUTION D'AUTOMOBILE. DE FAÇON ARBITRAIRE, L'UN DES RÉSULTATS EST APPELÉ **SUCCÈS** ET L'AUTRE **ÉCHEC**.



C'EST CE PROCESSUS ALÉATOIRE AUSSI APPELÉ ÉPREUVE DE BERNOULLI QUE NOUS ALLONS RÉPÉTER. L'EXPÉRIENCE QUI CONSISTE À RÉPÉTER LES ÉPREUVES EST APPELÉE

UN **schéma de Bernoulli**,
SI ELLE VÉRIFIE LES PROPRIÉTÉS
SUIVANTES :

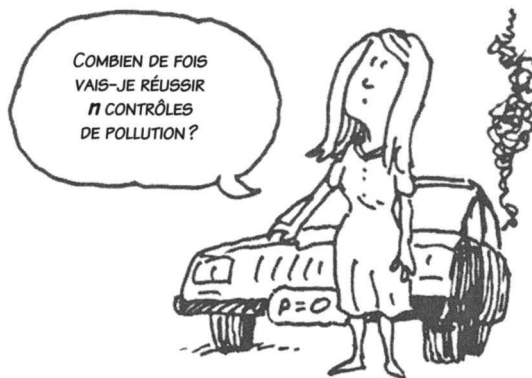
- 1) LE RÉSULTAT DE CHAQUE ÉPREUVE EST SOIT UN SUCCÈS, SOIT UN ÉCHEC.
- 2) LA PROBABILITÉ p DE SUCCÈS EST LA MÊME POUR CHAQUE ÉPREUVE.
- 3) LES ÉPREUVES SONT **INDÉPENDANTES** : LE RÉSULTAT D'UNE ÉPREUVE N'A AUCUNE INFLUENCE SUR LES AUTRES ÉPREUVES.



À PARTIR D'UNE ÉPREUVE DE BERNOULLI AYANT UN SUCCÈS AVEC PROBABILITÉ p , ON PEUT CONSTRUIRE UNE NOUVELLE VARIABLE ALÉATOIRE EN RÉPÉTANT LES ÉPREUVES.

La variable aléatoire binomiale X

EST LE NOMBRE DE SUCCÈS DANS UN SCHÉMA RÉPÉTÉ DE BERNOULLI À n ÉTAPES, OÙ p EST LA PROBABILITÉ DE SUCCÈS.

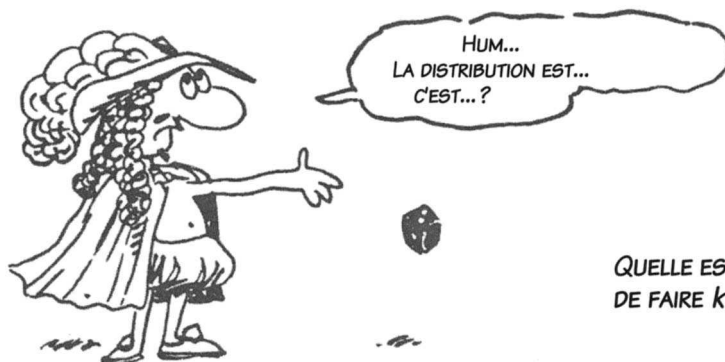


LE NOMBRE DE FACES (SUCCÈS) SUR DEUX LANCERS DE PIÈCE FOURNIT UN EXEMPLE DE VARIABLE ALÉATOIRE BINOMIALE. ICI $n = 2$ et $p = 0,5$.

$k = \text{NOMBRE DE SUCCÈS}$	0	1	2
$P(X = k)$	0,25	0,5	0,25



LE PREMIER PARI DE MÉRÉ FOURNIT UN AUTRE EXEMPLE. ON LANCE UN DÉ 4 FOIS DE SUITE. UN SUCCÈS CORRESPOND À FAIRE UN 6. LA DISTRIBUTION EST...



QUELLE EST LA PROBABILITÉ DE FAIRE k 6 EN 4 LANCERS ?

DE FAÇON GÉNÉRALE, QUELLE EST LA DISTRIBUTION DE PROBABILITÉS D'UNE VARIABLE ALÉATOIRE BINOMIALE X DE PROBABILITÉ p QUELCONQUE ET AVEC n ÉTAPES? UN CALCUL DE PROBABILITÉ FOURNIT LA RÉPONSE. LA PROBABILITÉ D'OBTENIR k SUCCÈS PARMI n ESSAIS, $P(X = k)$, EST :

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n - k}$$



ICI LE SYMBOLE $\binom{n}{k}$ SE LIT « COMBINAISON DE k PARMI n ». IL EST AUSSI APPELÉ **COEFFICIENT BINOMIAL**. EN FRANCE, ON LE NOTE SOUVENT C_n^k . IL CALCULE TOUTES LES DIFFÉRENTES FAÇONS D'OBTENIR k SUCCÈS SUR n ESSAIS. CHAQUE SÉQUENCE AVEC k SUCCÈS ET $n - k$ ÉCHECS A UNE PROBABILITÉ DE $p^k (1 - p)^{n - k}$ PAR LA LOI MULTIPLICATIVE. IL Y A $\binom{n}{k} = C_n^k$ DE CES SÉQUENCES.

AHHHHH...



LA FORMULE POUR LES COMBINAISONS EST :

$$C_n^k = \binom{n}{k} = \frac{n!}{k!(n - k)!}$$

où

$$n! = 1 \times 2 \times 3 \times \dots \times (n - 1) \times (n)$$

ET PAR CONVENTION $0! = 1$. PAR EXEMPLE, $\binom{4}{2}$ EST LE NOMBRE DE FAÇONS DE CHOISIR 2 LETTRES PARMI UN ENSEMBLE DE 4. CELA VAUT :

$$C_4^2 = \binom{4}{2} = \frac{4!}{2!2!} = \frac{24}{4} = 6$$

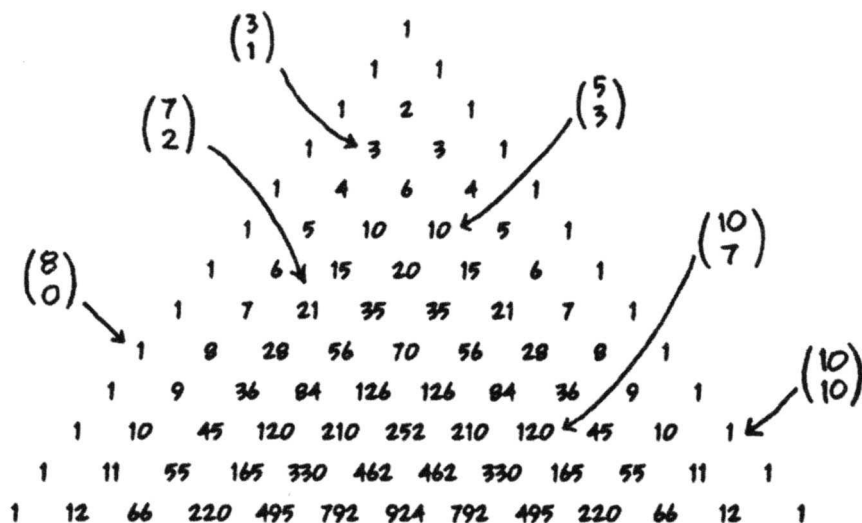
{A B C D}



AB AC AD

BC BD CD

ON PEUT AUSSI UTILISER LE **TRIANGLE DE PASCAL** POUR TROUVER LES COEFFICIENTS BINOMIAUX. CHAQUE ÉLÉMENT EST LA SOMME DES DEUX NOMBRES AU-DESSUS DE LUI.



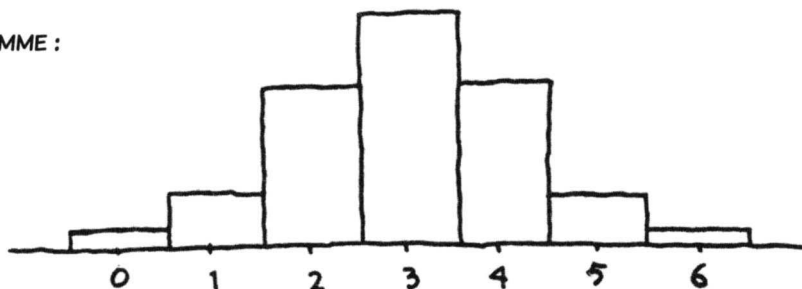
ETC.

POUR TROUVER $\binom{n}{k}$ ALLEZ À LA LIGNE n ET PRENEZ LA k^{e} VALEUR HORIZONTALEMENT (COMMENCEZ LE COMPTE TOUJOURS PAR ZÉRO).

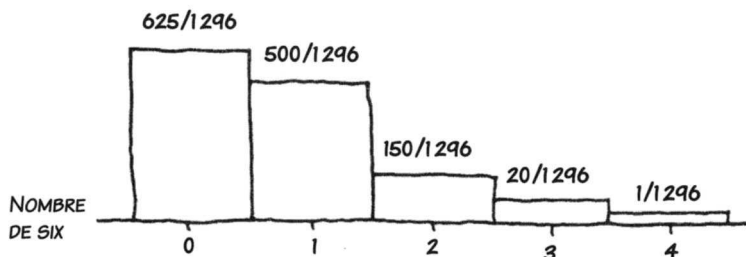
AVEC $p = 0,5$, LA DISTRIBUTION DE PROBABILITÉS BINOMIALE EST PARFAITEMENT SYMÉTRIQUE. AINSI, AVEC n , PAR EXEMPLE, ON OBTIENT :

$k = \text{NOMBRE DE FACES}$	0	1	2	3	4	5	6
$P(X = k)$	$\left(\frac{1}{2}\right)^6$	$6 \left(\frac{1}{2}\right)^6$	$15 \left(\frac{1}{2}\right)^6$	$20 \left(\frac{1}{2}\right)^6$	$15 \left(\frac{1}{2}\right)^6$	$6 \left(\frac{1}{2}\right)^6$	$\left(\frac{1}{2}\right)^6$

D'OÙ CET HISTOGRAMME :



POUR LE LANCER DES 4 DÉS DE MÉRÉ, LA DISTRIBUTION EST PLUS DÉSÉQUILIBRÉE.



LA MOYENNE ET LA VARIANCE DE LA LOI BINOMIALE SONT :

$$\mu = np$$

$$\sigma^2 = np(1 - p)$$

NOTEZ QUE LA MOYENNE EST INTUITIVE : AVEC n ÉPREUVES DE BERNOULLI, LE NOMBRE ESPÉRÉ DE SUCCÈS SERA np . LA VARIANCE PROVIENT DU FAIT QU'UN SCHÉMA BINOMIAL EST LA SOMME DE n ÉPREUVES DE BERNOULLI INDÉPENDANTES DE VARIANCE $p(1 - p)$.

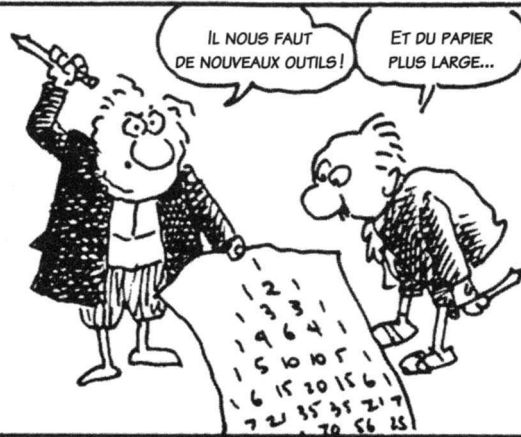


LES PARAMÈTRES D'UNE DISTRIBUTION BINOMIALE SONT n ET p . LA DISTRIBUTION, LA MOYENNE ET LA VARIANCE NE DÉPENDENT QUE DE CES DEUX NOMBRES. ON TROUVE LES TABLES DE LOI BINOMIALE DANS DES MANUELS ET DES PROGRAMMES INFORMATIQUES. VOICI LA TABLE POUR $n = 10$.

VALEUR DE $P(X = k)$

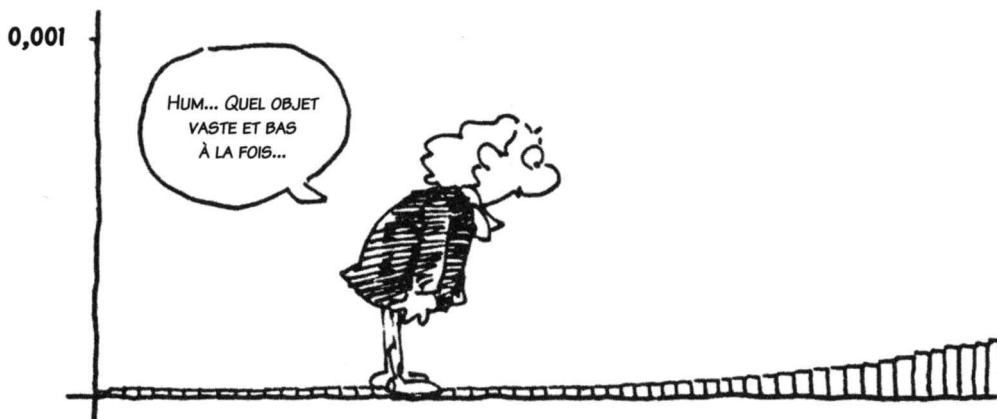
	k										
	0	1	2	3	4	5	6	7	8	9	10
P	0,1	0,344	0,387	0,194	0,057	0,011	0,001	0,000	0,000	0,000	0,000
	0,25	0,056	0,188	0,282	0,250	0,146	0,058	0,016	0,003	0,000	0,000
	0,5	0,001	0,010	0,044	0,117	0,205	0,246	0,205	0,117	0,044	0,010
	0,75	0,000	0,000	0,000	0,003	0,016	0,058	0,146	0,250	0,282	0,188
	0,9	0,000	0,000	0,000	0,000	0,000	0,001	0,011	0,057	0,194	0,344

MAIS FAIRE LE CALCUL
POUR DE GRANDES VALEURS
DE n S'AVÈRE TRÈS DÉLICAT...
DU MOINS AU XVIII^e SIÈCLE,
LORSQUE JACQUES BERNOULLI
(1654-1705) ET ABRAHAM
DE MOIVRE (1667-1754)
ESSAYAIENT DE LE FAIRE
SANS ORDINATEUR.



MOIVRE DÉPLOYA UNE NOUVELLE
ARME DE CALCUL ET MONTRA
QUE, LORSQUE $p = \frac{1}{2}$,
LA DISTRIBUTION BINOMIALE
POUVAIT ÊTRE APPROXIMÉE
PAR UNE FONCTION DE DENSITÉ
CONTINUE, TRÈS SIMPLE
À DÉCRIRE.

POUR VOIR COMMENT CELA FONCTIONNE, IMAGINEZ LA DISTRIBUTION BINOMIALE
AVEC $p = 0,5$ ET n TRÈS GRAND – DISONS 1 MILLION...



MAINTENANT, DIT MOIRE, FAITES GLISSER LE GRAPHE AFIN QUE SA MOYENNE SOIT ZÉRO.



ÉCRASEZ LA COURBE LE LONG DE L'AXE DES ABSCISSES AFIN QUE L'ÉCART-TYPE SOIT 1 TOUT EN L'ÉTIRANT SUR L'AXE DES ORDONNÉES POUR QUE L'AIRE SOIT ÉGALE À 1.



LE RÉSULTAT RESSEMBLE À UNE COURBE LISSE, SYMÉTRIQUE EN FORME DE CLOCHE, DONT MOIRE DONNA L'ÉQUATION SIMPLE :

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

CETTE FONCTION TRÈS IMPORTANTE EST APPELÉE LA **loi normale standard**.



UN OBJET
MAGNIFIQUE!



(e EST UNE CONSTANTE MATHÉMATIQUE UTILE QUI VAUT APPROXIMATIVEMENT 2,718).

POUR VÉRIFIER QUE LA COURBE EST BIEN EN FORME DE CLOCHE, NOTEZ QUE LOIN DE ZÉRO $f(z)$ EST QUASIMENT NULLE, PAR EXEMPLE $f(-5) = f(5) = 0,0000049$. LA COURBE EST SYMÉTRIQUE, CAR $f(z) = f(-z)$. SON MAXIMUM EST ATTEINT EN ZÉRO OÙ $f(0) = \frac{1}{\sqrt{2\pi}} = 0,39894$.

LA DISTRIBUTION EST APPELÉE LOI NORMALE **STANDARD** CAR LES DÉFORMATIONS ÉTAIENT ORGANISÉES POUR VÉRIFIER DES PROPRIÉTÉS SIMPLAS, QUE NOUS PRÉSENTONS SANS PREUVE :

$$\mu = 0$$

$$\sigma = 1$$

POUR RÉSUMER MOIVRE,
SI ON «**NORMALISE**»
LA DISTRIBUTION BINOMIALE AVEC
 $p = 1/2$ (AUTREMENT DIT, ON CENTRE
EN ZÉRO ET ON FAIT EN SORTE
QUE L'ÉCART-TYPE = 1),
ALORS ON OBTIENT PRESQUE
LA **DISTRIBUTION NORMALE**
STANDARD.

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$



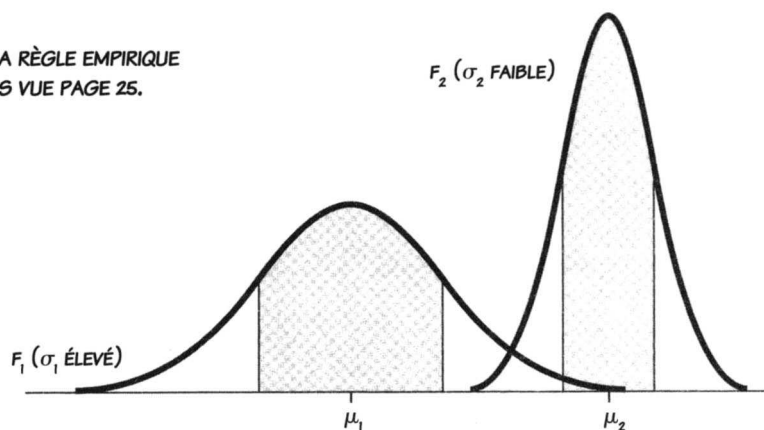
ON PEUT OBTENIR D'AUTRES LOIS NORMALES AVEC DES MOYENNES ET DES VARIANCES DIFFÉRENTES EN ÉTIRANT ET TRANSLATANT LA LOI NORMALE STANDARD. EN GÉNÉRAL, LA FORMULE SUIVANTE NOUS DONNE UNE DISTRIBUTION SYMÉTRIQUE, EN FORME DE CLOCHE CENTRÉE SUR LA MOYENNE μ ET D'ÉCART-TYPE σ :

$$f(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2}$$

VOICI DEUX DISTRIBUTIONS NORMALES DIFFÉRENTES, LES PARTIES GRISÉES REPRÉSENTENT LA RÉGION SITUÉE À MOINS D'UN ÉCART-TYPE DE LA MOYENNE. LEURS AIRES SONT ÉGALES. AINSI, POUR TOUTE VARIABLE ALÉATOIRE NORMALE, LA **PROBABILITÉ D'ÊTRE À MOINS D'UN ÉCART-TYPE DE LA MOYENNE EST TOUJOURS LA MÊME**, À SAVOIR À PEU PRÈS 0,68.

$$P(|X - \mu| < \sigma) \approx 0,68$$

CELA EXPLIQUE LA RÈGLE EMPIRIQUE QUE NOUS AVONS VUE PAGE 25.

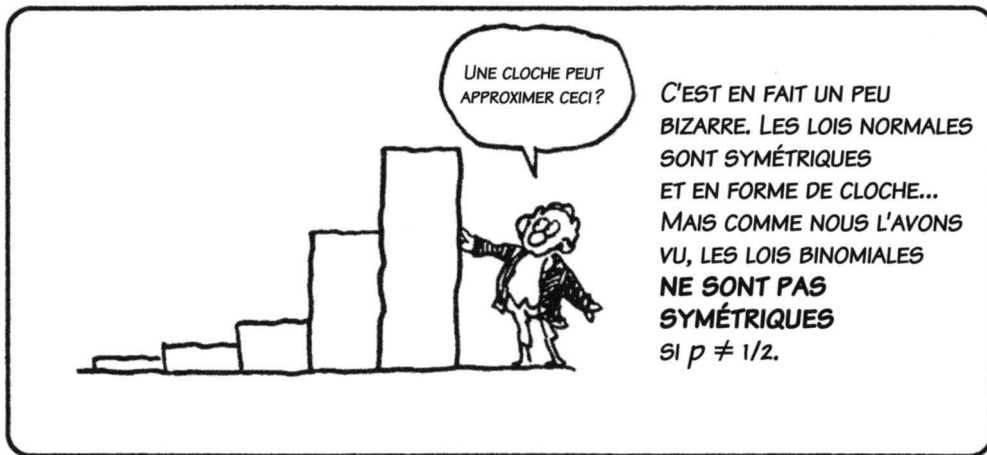


MOIVRE PROUVA QUE LA LOI NORMALE STANDARD AVAIT LA FORME DE LA LOI BINOMIALE (NORMALISÉE) POUR $p = 1/2$. MAIS EN FAIT, CELA FONCTIONNE POUR **TOUTE** VALEUR DE p .

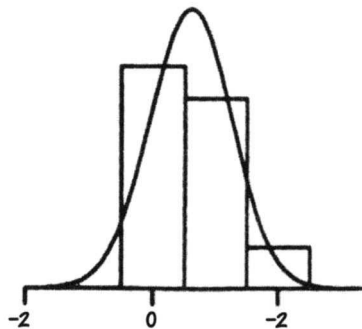
DE FAÇON GÉNÉRALE : QUEL QUE SOIT p , LA DISTRIBUTION BINOMIALE AVEC n ESSAIS DE PROBABILITÉ p EST APPROXIMÉE PAR UNE LOI NORMALE AVEC

$$\mu = np$$

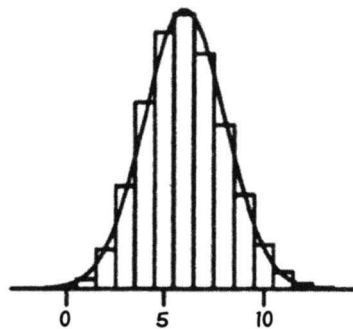
$$\text{ET } \sigma = \sqrt{np(1 - p)}.$$



MALGRÉ CELA, IL SE TROUVE QUE LORSQUE n AUGMENTE, L'ASYMÉTRIE DE LA LOI BINOMIALE EST CONTRARIÉE. COMME VOUS POUVEZ LE VOIR DANS CET EXEMPLE :



BINOMIALE : $n = 2$ ET $p = 0,3$



BINOMIALE : $n = 20$ ET $p = 0,3$

EN FAIT, LA DÉCOUVERTE FAITE PAR MOIVRE SUR LA LOI BINOMIALE EST UN CAS PARTICULIER D'UN RÉSULTAT ENCORE PLUS GÉNÉRAL QUI AIDE À MIEUX COMPRENDRE POURQUOI LA LOI NORMALE EST À LA FOIS SI IMPORTANTE ET SI RÉPANDUE DANS LA NATURE.

« Le théorème central limite » :

DES DONNÉES QUI DÉPENDENT DE NOMBREUX PETITS EFFETS ALÉATOIRES NON CORRÉLÉS SONT APPROXIMATIVEMENT DISTRIBUÉES NORMALEMENT.



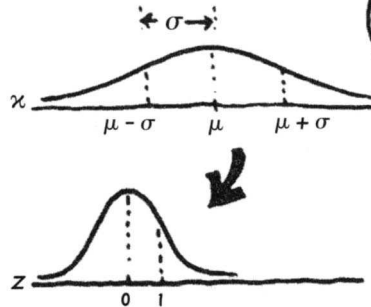
CELA EXPLIQUE QUE LA LOI NORMALE SE RETROUVE **PARTOUT** : DANS LES FLUCTUATIONS DES MARCHÉS FINANCIERS, LE POIDS DES ÉTUDIANTS, LES MOYENNES ANNUELLES DE TEMPÉRATURE, LES RÉSULTATS DU BAC : TOUTES CES QUANTITÉS SONT LE RÉSULTAT D'EFFETS NOMBREUX ET DIFFÉRENTS. PAR EXEMPLE, LE POIDS D'UN ÉTUDIANT EST LE RÉSULTAT DE SES GÈNES, DE SON ALIMENTATION, DE SES MALADIES ET DE SON PASSAGE AU PUB LA VEILLE. QUAND ON AJOUTE TOUS CES EFFETS, ON TROUVE UNE LOI NORMALE! (RAPPELEZ-VOUS QUE LA LOI BINOMIALE EST LE RÉSULTAT DE n ÉPREUVES INDÉPENDANTES DE BERNOULLI.)



LE Z-SCORE

$$z = \frac{x - \mu}{\sigma}$$

PERMET DE TRANSFORMER
UNE LOI NORMALE DE MOYENNE μ
ET D'ÉCART-TYPE σ EN UNE LOI
NORMALE CENTRÉE RÉDUITE,
DE MOYENNE 0 ET D'ÉCART-TYPE 1.



C'EST ENCORE
UNE OPÉRATION
DE GLISSEMENT
ET D'ÉTIREMENT...

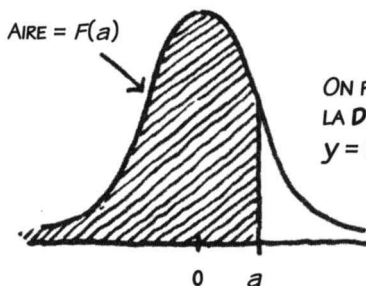


AINSI, POUR TROUVER LA PROBABILITÉ DE N'IMPORTE QUELLE LOI NORMALE, IL SUFFIT
DE CONNAÎTRE LA TABLE DE LA LOI NORMALE STANDARD $F(z)$.

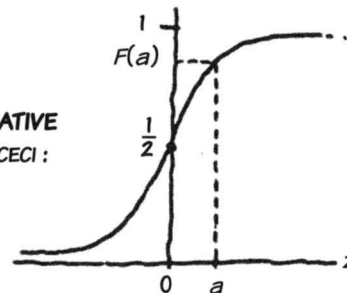
Z	-2,5	-2,4	-2,3	-2,2	-2,1	-2,0	-1,9	-1,8	-1,7	-1,6
F(z)	0,006	0,008	0,011	0,014	0,018	0,023	0,029	0,036	0,045	0,055
Z	-1,5	-1,4	-1,3	-1,2	-1,1	-1,0	-0,9	-0,8	-0,7	-0,6
F(z)	0,067	0,081	0,097	0,115	0,136	0,159	0,184	0,212	0,242	0,274
Z	-0,5	-0,4	-0,3	-0,2	-0,1	0,0	0,1	0,2	0,3	0,4
F(z)	0,309	0,345	0,382	0,421	0,460	0,500	0,540	0,579	0,618	0,655
Z	0,5	0,6	0,7	0,8	0,9	1,0	1,1	1,2	1,3	1,4
F(z)	0,691	0,726	0,758	0,788	0,816	0,841	0,864	0,885	0,903	0,919
Z	1,5	1,6	1,7	1,8	1,9	2,0	2,1	2,2	2,3	2,4
F(z)	0,933	0,945	0,955	0,964	0,971	0,977	0,982	0,986	0,989	0,992
Z	2,5									
F(z)	0,994									



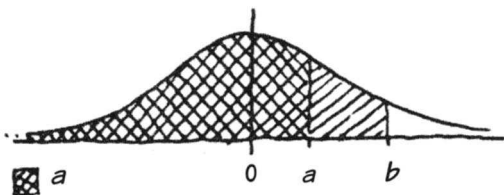
ICI $F(a) = P(z \leq a)$, SOIT L'AIRE COMPRISE ENTRE LA FONCTION
DE DENSITÉ ET L'AXE DES ABSCISSES ET À GAUCHE DE $z = a$.



ON PEUT AUSSI TRACER
LA DISTRIBUTION CUMULATIVE
 $y = F(z)$ QUI RESSEMBLE À CECI :



LA TABLE DE LOI NOUS PERMET DE CALCULER LA PROBABILITÉ QUE Z SOIT DANS L'INTERVALLE $a \leq z \leq b$. IL S'AGIT DE LA DIFFÉRENCE D'AIRES ENTRE $F(b)$ ET $F(a)$.



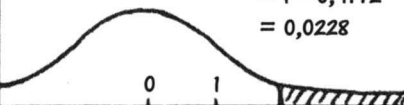
$$P(a \leq z \leq b) = F(b) - F(a)$$

AINSI, PAR EXEMPLE :

$$\begin{aligned} P(-1 < z < 1) &= F(1) - F(-1) \\ &= 0,8413 - 0,1587 \\ &= 0,6826 \end{aligned}$$



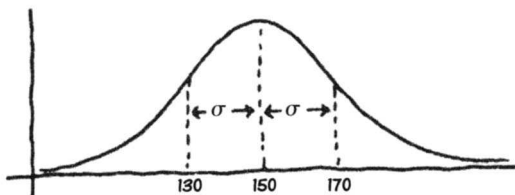
$$\begin{aligned} P(z > 2) &= 1 - F(2) \\ &= 1 - 0,9772 \\ &= 0,0228 \end{aligned}$$



ON PEUT UTILISER CETTE TABLE POUR TROUVER LES PROBABILITÉS DE N'IMPORTE QUELLE DISTRIBUTION NORMALE EN FAISANT LA TRANSFORMATION $z = (x - \mu)/\sigma$.



PAR EXEMPLE, SUPPOSONS QUE LE POIDS DES ÉTUDIANTS SOIT DISTRIBUÉ NORMALEMENT DE MOYENNE $\mu = 150$ ET D'ÉCART-TYPE $\sigma = 20$.

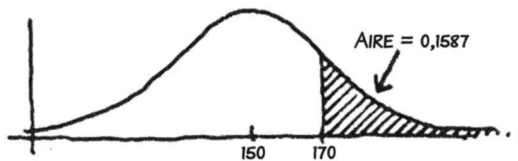


ALORS QUELLE EST LA PROBABILITÉ DE PESER PLUS DE 170 LIVRES ?

MAINTENANT, C'EST « SIMPLEMENT » DE L'ALGÈBRE.

$$\begin{aligned} P(X > 170) &= P\left(\frac{x - \mu}{\sigma} > \frac{170 - 150}{20}\right) \\ &= P\left(z > \frac{20}{20}\right) \\ &= P(z > 1) \end{aligned}$$

IL S'AGIT DE $1 - F(1)$. EN LISANT LA TABLE ON OBTIENT $1 - 0,8413 = 0,1587$:

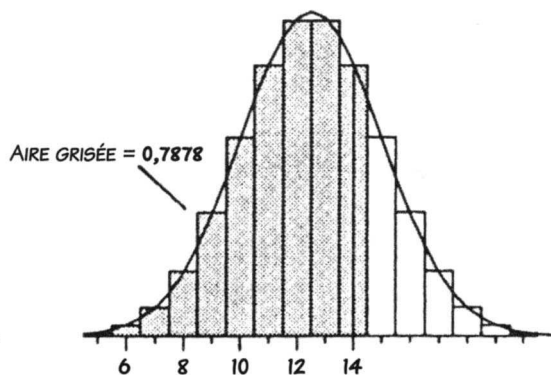


UN PEU MOINS D'UN ÉTUDIANT SUR SIX FAIT PENSER LA BALANCE AU-DESSUS DE 170 LIVRES (SOIT 77 kg).

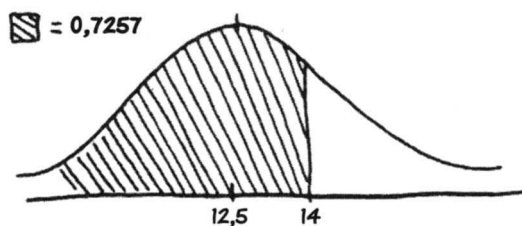
LA RÈGLE GÉNÉRALE POUR CALCULER DES PROBABILITÉS NORMALES EST DONC :

$$P(a \leq X \leq b) = F\left(\frac{b - \mu}{\sigma}\right) - F\left(\frac{a - \mu}{\sigma}\right)$$

MAINTENANT REVENONS
À L'APPROXIMATION BINOMIALE
DE MOIVRE. REGARDONS
LA DISTRIBUTION BINOMIALE
AVEC $p = \frac{1}{2}$ ET $n = 25$ (DISONS
25 TIRAGES DE PIÈCES).
ON PEUT CALCULER (OU LIRE
DANS UNE TABLE) N'IMPORTE
QUELLE PROBABILITÉ.
PAR EXEMPLE $P(X \leq 14)$
EST EXACTEMENT ÉGAL À **0,7878**.



CALCULONS MAINTENANT UNE VARIABLE ALÉATOIRE NORMALE X^* AVEC LA MÊME MOYENNE
 $\mu = np = (25)(0,5) = 12,5$ ET UN ÉCART-TYPE $\sigma = \sqrt{np(1-p)} = 2,5$.



$$\begin{aligned} P(X^* \leq 14) &= P\left(Z \leq \frac{14 - 12,5}{2,5}\right) \\ &= P(Z \leq 0,6) \\ &= \mathbf{0,7257} \end{aligned}$$



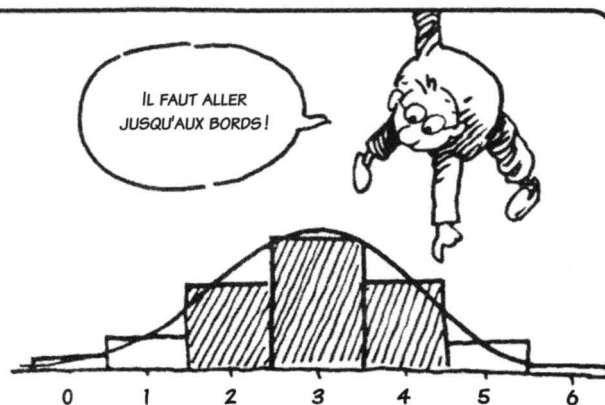
AH, MAIS ON PEUT FAIRE MIEUX !
SI ON REGARDE ATTENTIVEMENT
LE PREMIER HISTOGRAMME,
ON VOIT QUE LES BARRES SONT
CENTRÉES SUR LES ENTIERS.
CELA VEUT DIRE QUE $P(X^* \leq 14)$
EST EN FAIT L'AIRESITUÉE SOUS
LES BARRES EN DESSOUS
DE $x = 14,5$. NOUS DEVONS
DONC INCLURE CE 0,5. AINSI,

$$\begin{aligned} P(X^* \leq 14,5) &= P(Z \leq 0,8) \\ &= \mathbf{0,7881} \end{aligned}$$

UNE TRÈS BONNE APPROXIMATION
DE 0,7878, EN EFFET !

CE 0,5 SUPPLÉMENTAIRE
S'APPELLE LA **correction
de continuité**.

NOUS DEVONS L'INCLURE
POUR OBTENIR UNE BONNE
APPROXIMATION CONTINUE
DE NOTRE LOI BINOMIALE
DISCRÈTE X. CELA SE RÉSUME
AVEC UNE ÉQUATION
UN PEU HORRIBLE
MAIS SIMPLE D'UTILISATION.

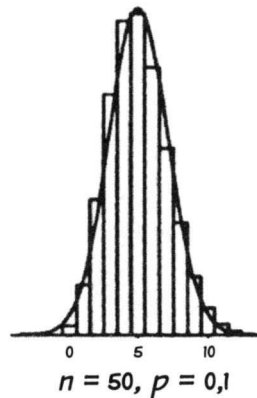
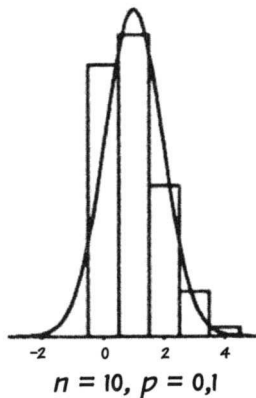
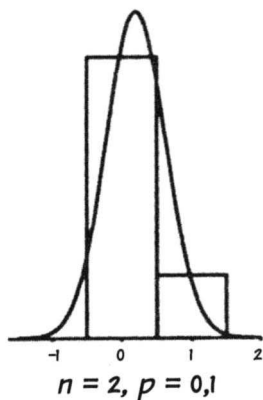


$$P(a \leq X \leq b) \approx P\left(\frac{a - \frac{1}{2} - np}{\sqrt{np(1-p)}} \leq Z \leq \frac{b + \frac{1}{2} - np}{\sqrt{np(1-p)}}\right)$$

QUAND CETTE APPROXIMATION EST-ELLE « SUFFISAMMENT BONNE » ? POUR LES STATISTICIENS LA RÈGLE EMPIRIQUE VEUT QUE n SOIT SUFFISAMMENT GRAND POUR QUE LES NOMBRES ESPÉRÉS DES SUCCÈS ET DES ÉCHECS SOIENT TOUS DEUX **SUPÉRIEURS À 5** :

$$np \geq 5 \quad \text{ET} \quad n(1-p) \geq 5$$

LES HISTOGRAMMES SUIVANTS ILLUSTRENT QUE, LORSQUE $p = 0,1$, L'APPROXIMATION EST MÉDIOCRE (VOIRE PIRE) JUSQU'À CE QUE n ATTEIGNE 50, POUR LEQUEL $np = 5$.



POURQUOI L'APPROXIMATION BINOMIALE EST-ELLE SI IMPORTANTE? EN FAIT LA DISTRIBUTION BINOMIALE APPARAÎT SOUVENT DANS LA NATURE, CE QUI N'EST PAS SI COMPLIQUÉ À COMPRENDRE. MAIS ELLE PEUT ÊTRE FASTIDIEUSE À CALCULER.



IL Y EN A UNE NOUVELLE
POUR CHAQUE VALEUR
DE n ET DE p ...

LA LOI NORMALE QUI L'APPROXIME EST PEUT-ÊTRE MOINS INTUITIVE, MAIS ELLE EST TRÈS SIMPLE D'UTILISATION. GRÂCE À LA TRANSFORMATION EN Z, ON PEUT CONVERTIR TOUTES LES LOIS NORMALES EN LOI NORMALE STANDARD, ET ON PEUT DONC TROUVER LES PROBABILITÉS ASSOCIÉES DIRECTEMENT À PARTIR D'UNE SEULE TABLE NUMÉRIQUE.

DANS UN LIVRE
OU SUR L'ÉCRAN
D'UN ORDINATEUR!



EN OUTRE, LA LOI NORMALE EST VRAIMENT
LA MÈRE DE TOUTES LES DISTRIBUTIONS !

C'EST LE THÉORÈME
CENTRAL LIMITE!

MAMAN!
MAMAN!



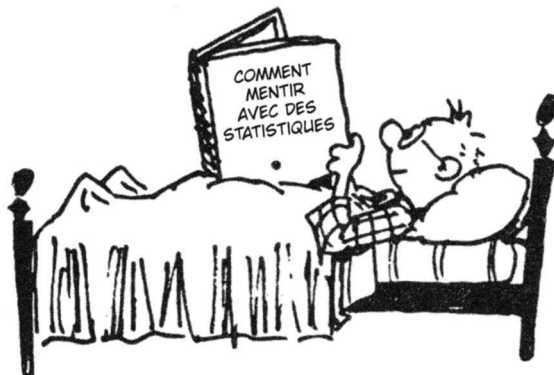
Chapitre 6

Échantillonnage

À PRÉSENT, APRÈS UN RÉGIME VARIÉ DE PIÈCES, DE DÉS ET D'IDÉES ABSTRAITES VOUS DEVEZ VOUS DEMANDER EN QUOI LES OUTILS STATISTIQUES QUE NOUS AVONS CONSTRUITS PEUVENT NOUS AIDER DANS LE **MONDE RÉEL**. EH BIEN, NOUS ALLONS ENFIN LE DÉCOUVRIR...

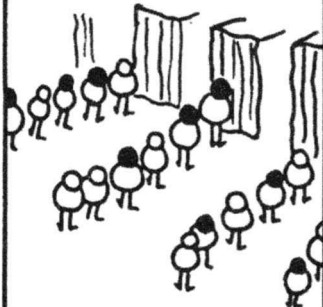


DANS CE CHAPITRE, NOUS COMMENÇONS À REGARDER LE **VRAI CŒUR** DU BUSINESS DES STATISTIQUES, DONT LE BUT, APRÈS TOUT, EST DE FAIRE GAGNER DU **TEMPS** ET DE L'**ARGENT**. LES GENS DÉTESTENT PERDRE LEUR TEMPS EN FAISANT UN **TRAVAIL INUTILE**. S'IL Y A UNE CHOSE DONT LES STATISTIQUES SONT CAPABLES, C'EST JUSTEMENT DE NOUS INDiquer JUSQU'À QUEL POINT NOUS POUVONS NOUS PERMETTRE D'ÊTRE PARESSEUX.

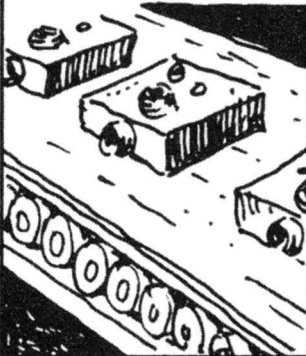


LE PROBLÈME AVEC NOTRE MONDE EST QUE LES COLLECTIONS DE CHOSES SONT TELLEMENT VASTES QU'IL EST DIFFICILE D'OBTENIR L'INFORMATION DONT ON A BESOIN.

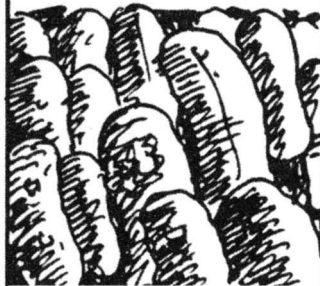
POPULATION DES VOTANTS :
QUEL POURCENTAGE POUR
QUEL CANDIDAT ?



PRODUITS MANUFACTURÉS :
QUELLE PROPORTION
DE PRODUITS DÉFECTUEUX ?



CORNICHONS : QUELLE EST
LEUR LONGUEUR MOYENNE ?



L'INFORMATION EST UTILE
AUX FABRICANTS DE BOCAUX
DE CORNICHONS !

LA RÉPONSE D'UN CASTOR
APPLIQUÉ, TRAVAILLEUR
ET CANDIDE SERAIT
DE MESURER CHAQUE
CORNICHON DU MONDE
ET DE FAIRE UN PEU
D'ARITHMÉTIQUE.

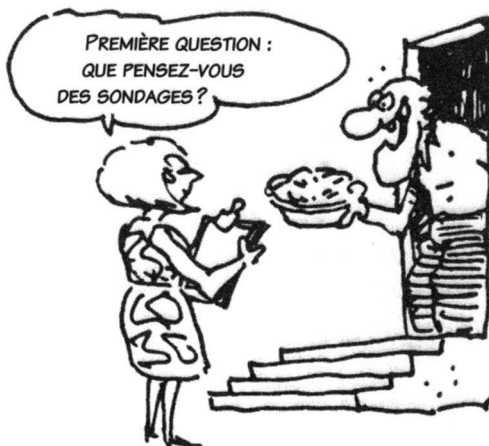


MAIS NOUS NE SOMMES PAS DES CASTORS
— NOUS SOMMES DES STATISTICIENS !
NOUS CHERCHONS UN MOYEN PLUS FACILE.



OH, EH BIEN,
J'AI MANGÉ LE STYLO
DE TOUTE FAÇON...

UNE MÉTHODE EST DE PRENDRE
UN **ÉCHANTILLON**...
UN SOUS-ENSEMBLE RELATIVEMENT PETIT
DE LA POPULATION TOTALE, À LA FAÇON
DE CE QUE FONT LES SONDEURS
POUR LES ÉLECTIONS.



UNE QUESTION ÉVIDENTE APPARAÎT : QUELLE DOIT ÊTRE LA TAILLE DE L'ÉCHANTILLON POUR AVOIR DES RÉSULTATS
SIGNIFICATIFS ?



ET LA RÉPONSE À CETTE QUESTION,
QUE VOUS DEVRIEZ GRAVER
DANS VOTRE CERVEAU, EST : SI n
EST LE NOMBRE D'ÉLÉMENTS
DE L'ÉCHANTILLON, ALORS TOUT
DÉPEND DE LA VALEUR :

$$\frac{1}{\sqrt{n}}$$

$\frac{1}{\sqrt{n}}$?
JE NE SAVAIS MÊME
PAS QU'IL PARTICIPAIT
AU SCRUTIN !



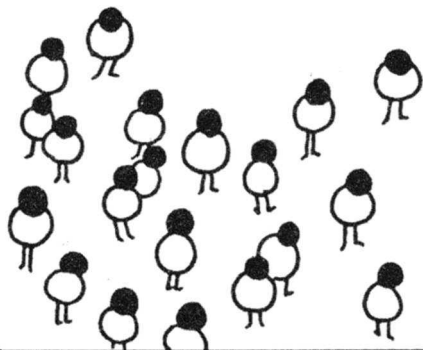
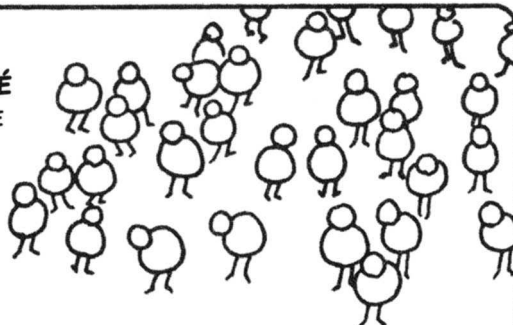
MÉTHODE D'ÉCHANTILLONNAGE



ENCORE
PLUS UTILE
QUE LE TRICOT!

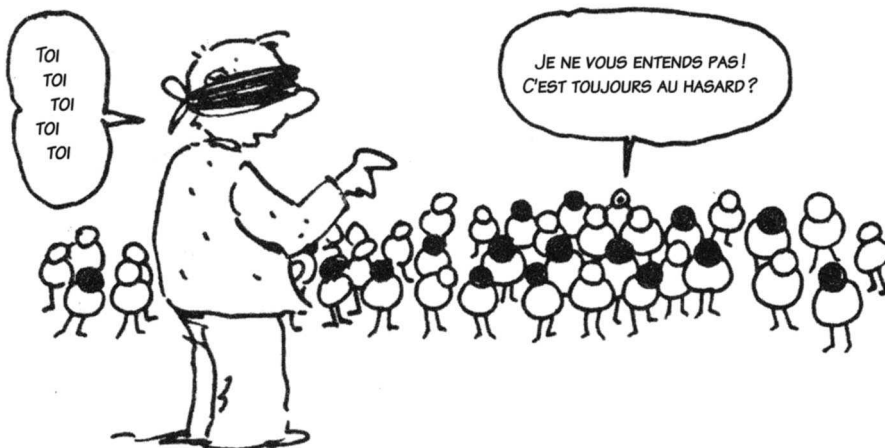
AVANT D'EN VENIR AUX NOMBRES,
NOUS DEVONS SIGNALER QUE LA **QUALITÉ**
DE L'ÉCHANTILLON EST AUSSI IMPORTANTE
QUE SA **TAILLE**.

COMMENT S'ASSURER QUE NOUS
CHOISSONS UN ÉCHANTILLON
REPRÉSENTATIF?



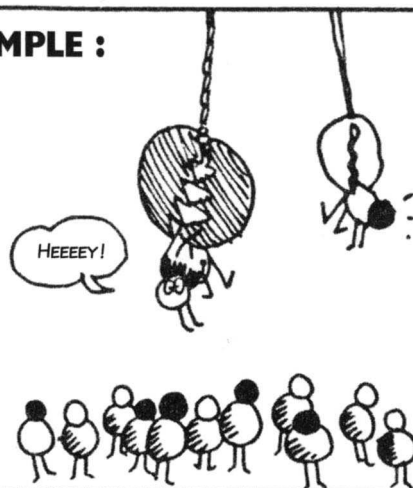
LE PROCESSUS DE **SÉLECTION** EST
LUI-MÊME ESSENTIEL. PAR EXEMPLE,
UN SONDAGE ÉLECTORAL QUI EXCLURAIT
SYSTÉMATIQUEMENT LES PARISIENS SERAIT
INUTILE. IL Y A BIEN D'AUTRES MANIÈRES
DE GÂCHER OU DE BIAISER UN ÉCHANTILLON.

NE PROLONGEONS PAS LE MYSTÈRE : LA FAÇON D'OBTENIR DES RÉSULTATS STATISTIQUES
FIABLES EST DE PRENDRE UN ÉCHANTILLON AU **hasard**.



Échantillonnage aléatoire SIMPLE :

IL NÉCESSITE UNE GRANDE POPULATION D'ÉLÉMENTS ET UNE PROCÉDURE POUR EN CHOISIR n PARMI EUX. SI LA PROCÉDURE ASSURE QUE CHACUN DES ÉCHANTILLONS POSSIBLES DE n ÉLÉMENTS EST ÉQUIPROBABLE, ALORS ELLE NOUS PROCURE UN **échantillon aléatoire simple**.



L'ÉCHANTILLONNAGE ALÉATOIRE SIMPLE DÉFINIT UN STANDARD PAR RAPPORT AUX AUTRES MÉTHODES, CAR IL VÉRIFIE DEUX PROPRIÉTÉS :



1) NON BIAISÉ : CHAQUE ÉLÉMENT A LES MÊMES CHANCES D'ÊTRE SÉLECTIONNÉ.

2) INDÉPENDANCE : LA SÉLECTION D'UN ÉLÉMENT N'A AUCUNE INCIDENCE SUR LA SÉLECTION DES AUTRES ÉLÉMENTS.

MALHEUREUSEMENT, IL EST DIFFICILE DE TROUVER DES ÉCHANTILLONS INDÉPENDANTS ET TOTALEMENT SANS BIAIS DANS LE MONDE RÉEL. PAR EXEMPLE, UN SONDAGE EFFECTUÉ EN COMPOSANT AU HASARD DES NUMÉROS DE TÉLÉPHONE SERA BIAISÉ : IL IGNORE LES PERSONNES NE DÉTENANT PAS DE TÉLÉPHONE ET SURREPRÉSENTE CELLES AYANT PLUS D'UNE LIGNE TÉLÉPHONIQUE.



THÉORIQUEMENT, IL EST POSSIBLE D'OBTENIR UN ÉCHANTILLON ALÉATOIRE SIMPLE EN ÉTABLISSANT LA **BASE D'ÉCHANTILLONNAGE** : UNE LISTE DE TOUS LES ÉLÉMENTS DE LA POPULATION. ON PEUT ALORS AVEC UN GÉNÉRATEUR DE NOMBRES ALÉATOIRES SÉLECTIONNER n ÉLÉMENTS AU HASARD.



DE MÊME, ON PEUT INSCRIRE LES NOMS SUR DES CARTES ET EN PRÉLEVER n AU HASARD DANS UNE URNE OU UN TAMBOUR.

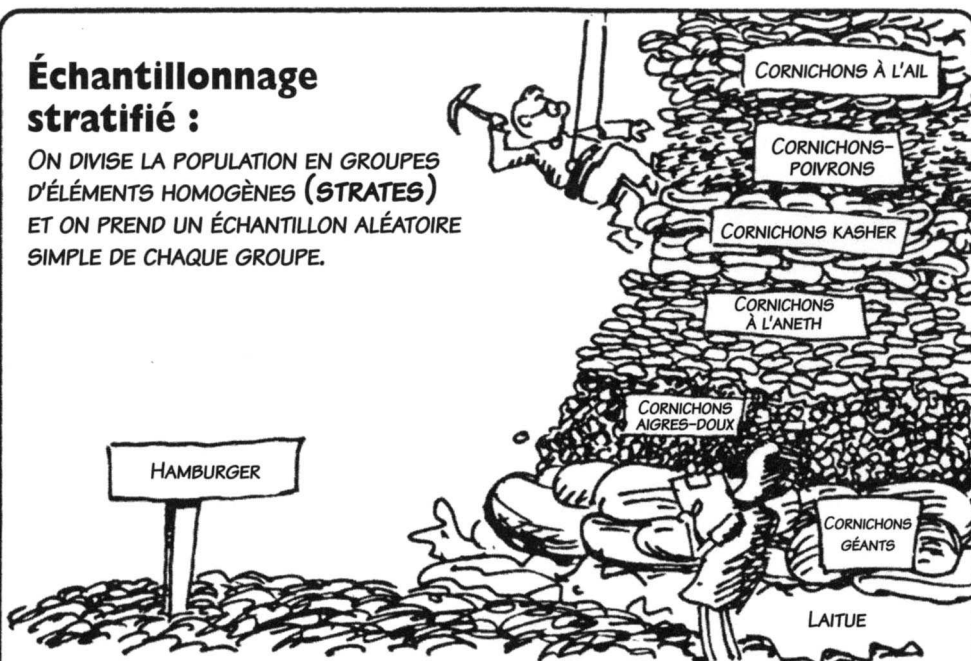
MAIS CE N'EST PAS TOUJOURS FACILE. IL PEUT ÊTRE TROP COÛTEUX, POLÉMIQUE OU MÊME IMPOSSIBLE DE CRÉER LA BASE D'ÉCHANTILLONNAGE. PAR EXEMPLE, UNE ÉTUDE MINISTÉRIELLE SUR LA QUALITÉ DES EAUX NÉCESSITE UNE BASE D'ÉCHANTILLONNAGE DES LACS. IL FAUT DONC QUE QUELQU'UN DÉCIDE :



Y A-T-IL D'AUTRES FAÇONS PLUS EFFICACES ET MOINS COÛTEUSES DE CRÉER UN ÉCHANTILLON ? LA RÉPONSE EST OUI SI VOUS CONNAISSEZ UN PEU VOTRE POPULATION. PAR EXEMPLE...

Échantillonnage stratifié :

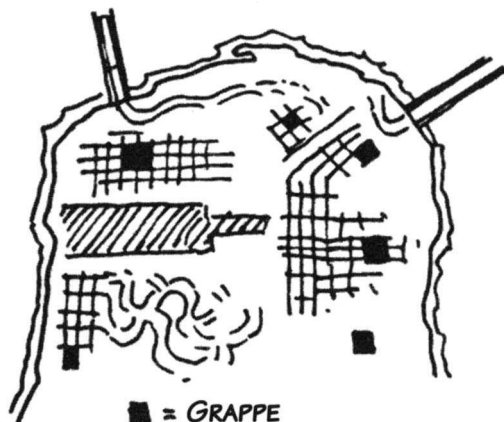
ON DIVISE LA POPULATION EN GROUPES D'ÉLÉMENTS HOMOGÈNES (**STRATES**) ET ON PREND UN ÉCHANTILLON ALÉATOIRE SIMPLE DE CHAQUE GROUPE.



PAR EXEMPLE, LA POPULATION DE TOUS LES **CORNICHONS** PEUT ÊTRE STRATIFIÉE PAR **TYPES**. LEUR TAILLE SERA ALORS MOINS VARIABLE À L'INTÉRIEUR DE CHAQUE STRATE.

Échantillonnage en grappes :

ON SUBDIVISE LA POPULATION EN PLUS PETITES GRAPPES. ON PREND ALORS UN ÉCHANTILLON ALÉATOIRE SIMPLE DE GRAPPES ET ON OBSERVE TOUS LES ÉLÉMENTS DES GRAPPES SÉLECTIONNÉES. CETTE TECHNIQUE EST TRÈS RENTABLE LORSQUE LES COÛTS DE TRANSPORT ENTRE LES ÉLÉMENTS SONT ÉLEVÉS.

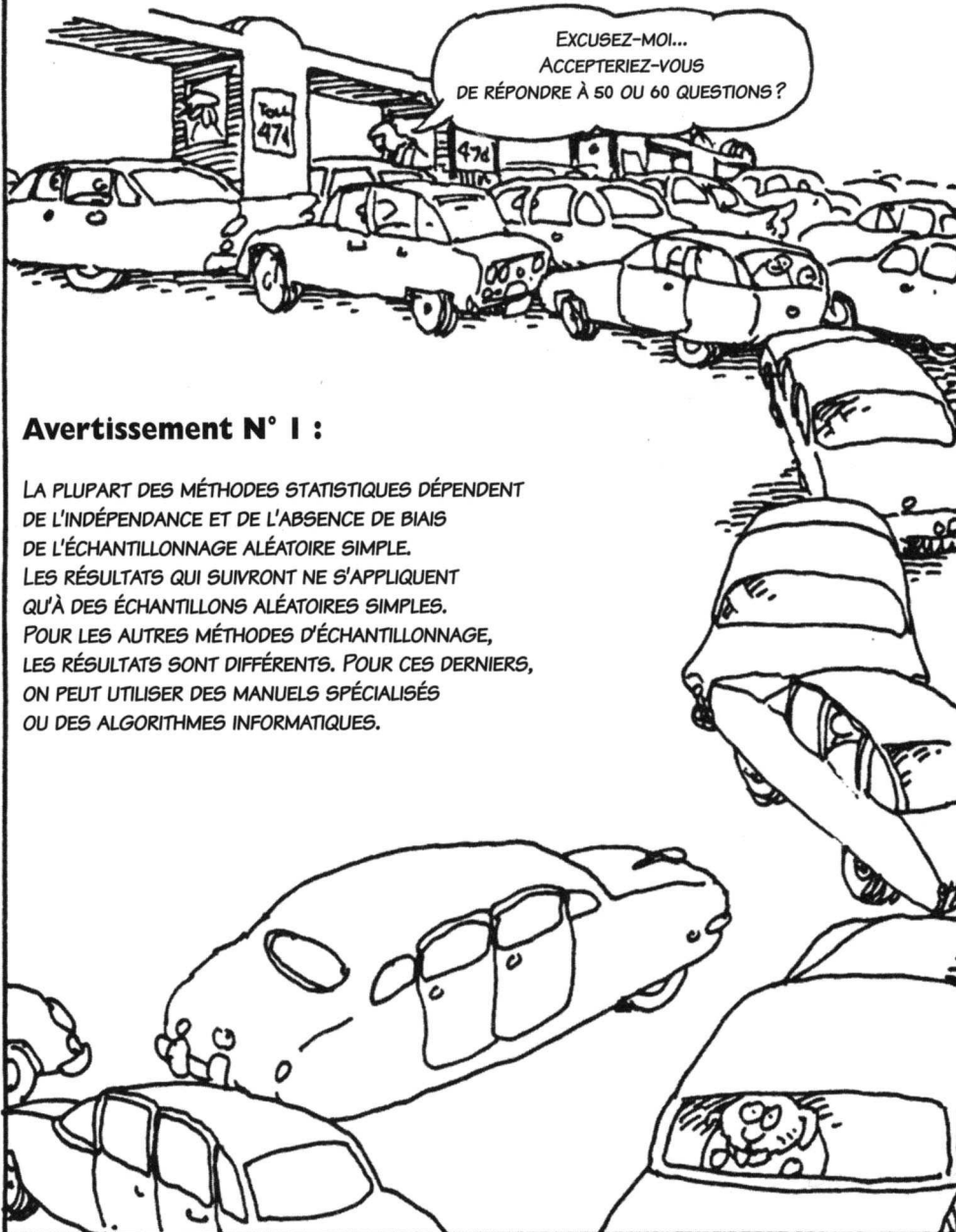


POUR UNE ÉTUDE SUR LES FOYERS D'UNE VILLE, ON PEUT PAR EXEMPLE : DIVISER LA VILLE EN BLOCS, SÉLECTIONNER UN ÉCHANTILLON DE BLOCS ET SONDER CHAQUE FOYER DES BLOCS CHOISIS.

Échantillonnage systématique :

ON COMMENCE PAR CHOISIR UN PREMIER ÉLÉMENT AU HASARD, PUIS ON PREND TOUS LES k^e SUIVANTS. PAR EXEMPLE, UNE ÉTUDE SUR LE TRAFIC AUTOROUTIER

PEUT TESTER CHAQUE VOITURE SUR 100 PASSAGES À UNE BARRIÈRE DE PÉAGE. CETTE STRATÉGIE EST FACILE À METTRE EN PLACE ET ELLE EST PLUS EFFICACE SI LA CIRCULATION VARIE PEU AVEC LE TEMPS.



Avertissement N° 1 :

LA PLUPART DES MÉTHODES STATISTIQUES DÉPENDENT DE L'INDÉPENDANCE ET DE L'ABSENCE DE BIAIS DE L'ÉCHANTILLONNAGE ALÉATOIRE SIMPLE. LES RÉSULTATS QUI SUIVRONT NE S'APPLIQUENT QU'À DES ÉCHANTILLONS ALÉATOIRES SIMPLES. POUR LES AUTRES MÉTHODES D'ÉCHANTILLONNAGE, LES RÉSULTATS SONT DIFFÉRENTS. POUR CES DERNIERS, ON PEUT UTILISER DES MANUELS SPÉCIALISÉS OU DES ALGORITHMES INFORMATIQUES.

Avertissement N° 2 :



SANS CRÉATION DE HASARD,
IL N'Y A PAS D'ANALYSE STATISTIQUE
FIABLE, ET CE QUELLES QUE SOIENT
LES ADAPTATIONS. L'AVANTAGE
DE L'ÉCHANTILLONNAGE ALÉATOIRE
EST QUE CELUI-CI GARANTIT
STATISTIQUEMENT LA RIGUEUR
DES ENQUÊTES.

UNE MÉTHODE TRÈS SOUVENT UTILISÉE ET PARTICULIÈREMENT SUJETTE À BIAIS EST CELLE
DE L'**échantillonnage de commodité**. L'ENQUÊTEUR ÉVITE LES TRACAS
D'UNE PROCÉDURE PARTICULIÈRE EN PRENANT SIMPLEMENT
LES 11 PREMIERS ÉLÉMENTS DE LA POPULATION
QUI LUI TOMBENT SOUS LA MAIN.

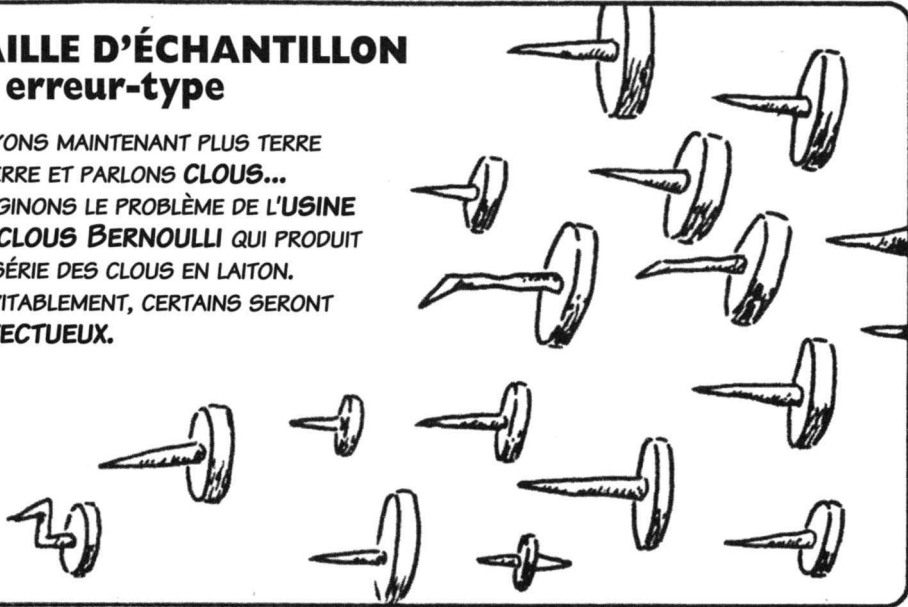


ON EN TROUVE UN EXEMPLE DANS LE LIVRE *LES FEMMES ET L'AMOUR* ÉCRIT EN 1987 PAR SHERE HITE.
EN EFFET, 100 000 QUESTIONNAIRES FURENT ENVOYÉS À DES ORGANISATIONS FÉMININES
(UN **ÉCHANTILLON D'OPPORTUNITÉ**). SEULEMENT 4,5 % FURENT RENVOYÉS (**BIAIS DANS
LES RÉPONSES**). SES « RÉSULTATS » ÉTAIENT DONC FONDÉS SUR UN ÉCHANTILLON DE FEMMES
MOTIVÉES, DÉSIREUSES DE RÉPONDRE AUX QUESTIONS DE L'ENQUÊTE POUR DES RAISONS DIVERSES.

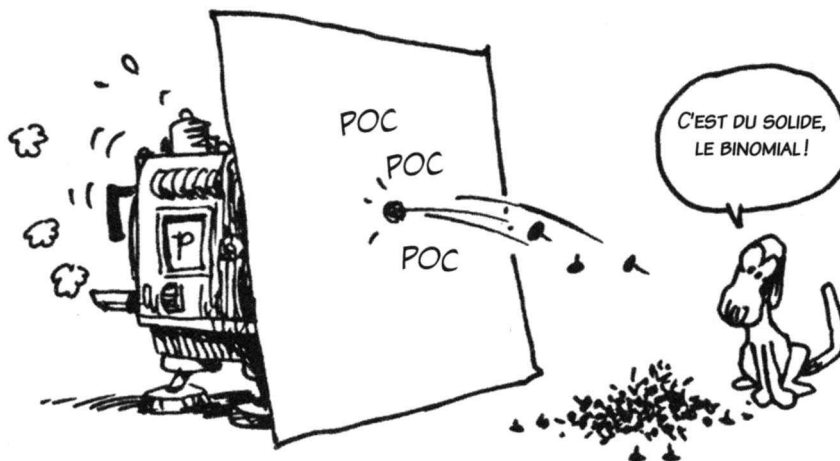


TAILLE D'ÉCHANTILLON et erreur-type

SOYONS MAINTENANT PLUS TERRE
À TERRE ET PARLONS **CLOUS...**
IMAGINONS LE PROBLÈME DE L'**USINE
DE CLOUS BERNOULLI** QUI PRODUIT
EN SÉRIE DES CLOUS EN LAITON.
INÉVITABLEMENT, CERTAINS SERONT
DÉFECTUEUX.



LES LECTEURS ATTENTIFS RECONNAÎTRONT UN **SCHEMA DE BERNOULLI** :
CHAQUE NOUVEAU CLOU EST LE RÉSULTAT D'UNE ÉPREUVE DE BERNOULLI
AVEC UNE PROBABILITÉ p DE SUCCÈS (CLOU SANS DÉFAUT) ET UNE PROBABILITÉ
 $1 - p$ D'ÉCHEC (CLOU DÉFECTUEUX).



ON PEUT CONSIDÉRER QUE TOUT SE PASSE COMME SI UNE **MACHINE DE BERNOULLI
CACHÉE, MAIS RÉELLE**, PRODUISAIT SELON UNE PROBABILITÉ p LES RÉSULTATS
OBSERVÉS DANS LE MONDE RÉEL.

COMME LA MACHINE DE BERNOULLI
EST INVISIBLE, ON NE CONNAÎT PAS
 p MAIS ON AIMERAIT LE CONNAÎTRE.
ON PREND DONC UN ÉCHANTILLON
ALÉATOIRE DE n CLOUS,
ET ON REMARQUE QUE x D'ENTRE
EUX SONT IRRÉPROCHABLES.



HUM...
JE CROIS QUE $n = 400$
ET $x = 352$...

MAINTENANT LA PROPORTION DE SUCCÈS DANS L'ÉCHANTILLON DEVRAIT ÊTRE
DE L'ORDRE DE p . ON NOTE CETTE PROPORTION \hat{p} ET ON PRONONCE « p -CHAPEAU ».

$$\hat{p} = \frac{x}{n}$$

\hat{p} EST LE RATIO DU NOMBRE x DE SUCCÈS DE L'ÉCHANTILLON PAR LA TAILLE n
DE L'ÉCHANTILLON. PAR EXEMPLE, SI p VAUT 0,85 ET QUE L'ON AIT ÉCHANTILLONNÉ
 $n = 1000$ CLOUS, ON POURRAIT TROUVER $x = 832$ CLOUS SANS DÉFAUT, DE SORTE
QUE $\hat{p} = 0,832$.

LA QUESTION EST :
CETTE ESTIMATION
EST-ELLE BONNE ?



AÏE!
EST-CE VRAIMENT
« BON » ?

ET NOUS ALLONS RÉPONDRE
À UNE AUTRE QUESTION :
QUE SIGNIFIE NOTRE PREMIÈRE
QUESTION ?

ON NE PEUT PAS CONNAÎTRE LA VÉRITABLE DIFFÉRENCE ENTRE \hat{p} ET p , CAR NOUS NE SAVONS PAS LA VALEUR EXACTE DE p . LA VRAIE QUESTION EST : SI L'ON PREND **PLUSIEURS ÉCHANTILLONS DE 1000 CLOUS** ET QUE L'ON OBSERVE CHAQUE \hat{p} , COMMENT CES \hat{p} SERONT-ILS DISTRIBUÉS AUTOUR DE p ?



EN FAIT, CES VALEURS DE \hat{p} RESSEMBLENT FORTEMENT À UNE **VARIABLE ALÉATOIRE** : LA SÉLECTION D'ÉCHANTILLON DE n ÉLÉMENTS EST UNE EXPÉRIENCE ALÉATOIRE ET L'OBSERVATION DE \hat{p} EN EST LE RÉSULTAT NUMÉRIQUE !



POUR ÊTRE PRÉCIS, SI X EST LE NOMBRE DE SUCCÈS DE L'ÉCHANTILLON, ALORS X N'EST RIEN D'AUTRE QU'UNE VARIABLE ALÉATOIRE BINOMIALE (AVEC n ESSAIS ET DE PROBABILITÉ p)... ET NOUS DÉFINISSONS LA **PROPORTION D'ÉCHANTILLON** (OU PROPORTION OBSERVÉE) COMME ÉTANT LA VARIABLE ALÉATOIRE :

$$\hat{p} = \frac{X}{n}$$

GRAND \hat{p} EST LA VARIABLE ALÉATOIRE, PETIT \hat{p} EST SA VALEUR POUR UN ÉCHANTILLON DONNÉ !



CONNAISSANT X , ON CONCLUT RAPIDEMENT QUELQUES PROPRIÉTÉS SUR \hat{P} :

1) LA MOYENNE DE \hat{P} EST $E[\hat{P}] = p$

2) L'ÉCART-TYPE DE \hat{P} EST

$$\sigma(\hat{P}) = \frac{\sqrt{p(1-p)}}{\sqrt{n}}$$

3) POUR DE GRANDES VALEURS DE n ,
 \hat{P} EST APPROXIMATIVEMENT NORMALE



ET VOILÀ VOUS SAVEZ TOUT! LES VALEURS OBSERVÉES DE \hat{P} SONT CENTRÉES EN p (SANS SURPRISE), ET L'ÉCART-TYPE (OU EN ANGLAIS SPREAD) EST PROPORTIONNEL À CE NOMBRE MAGIQUE QUE NOUS AVONS MENTIONNÉ EN DÉBUT DE CHAPITRE.



$$\frac{1}{\sqrt{n}}$$



ET COMME \hat{P} EST PRESQUE NORMALE, ON PEUT UTILISER LA RÈGLE EMPIRIQUE QUI DIT QU'APPROXIMATIVEMENT 68 % DE NOS ESTIMATIONS D'ÉCHANTILLON SERONT À UN ÉCART-TYPE DE LA VRAIE VALEUR p .

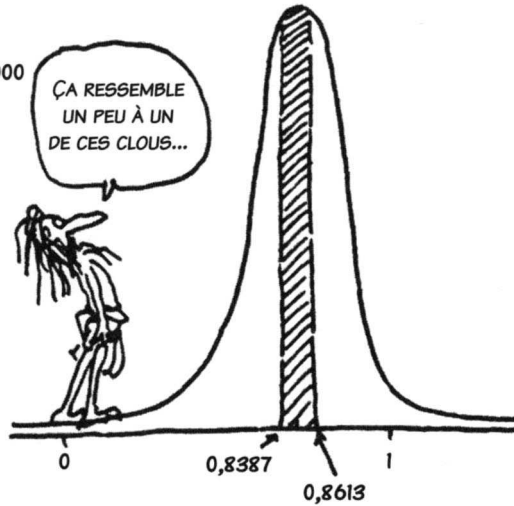


REVENONS À NOS CLOUS AVEC $n = 1000$
ET $p = 0,85$, L'ÉCART-TYPE EST DE :

$$\sigma(\hat{p}) = \sqrt{\frac{(0,85)(0,15)}{1000}}$$

$$= 0,0113$$

68 % DE NOS ESTIMATIONS
PONCTUELLES DOIVENT
ÊTRE DANS L'INTERVALLE
 $0,8387 \leq \hat{p} \leq 0,8613$



L'ÉCART-TYPE DE \hat{p} EST UNE MESURE
DE L'**erreur-type d'échantillonnage**.
COMME NOUS L'AVONS VU, CETTE ERREUR
D'ÉCHANTILLONNAGE EST INVERSEMENT PROPORTIONNELLE
À \sqrt{n} . AUGMENTER LA TAILLE DE L'ÉCHANTILLON
D'UN FACTEUR 4 RÉDUIT L'ERREUR-TYPE $\sigma(\hat{p})$ DE MOITIÉ.

DÉJÀ, AVEC $n = 100$,
ON VOIT QUE $\sigma(\hat{p})$ EST
DE L'ORDRE DE $3\frac{1}{2}\%$.

TAILLE D'ÉCHANTILLON DES CLOUS, $p = 0,85$

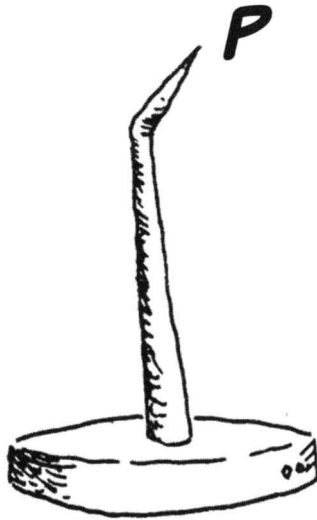
n	1	4	16	25	100	10 000
\sqrt{n}	1	2	4	5	10	100
$\sigma(\hat{p})$	0,357	0,1785	0,089	0,071	0,0357	0,0036



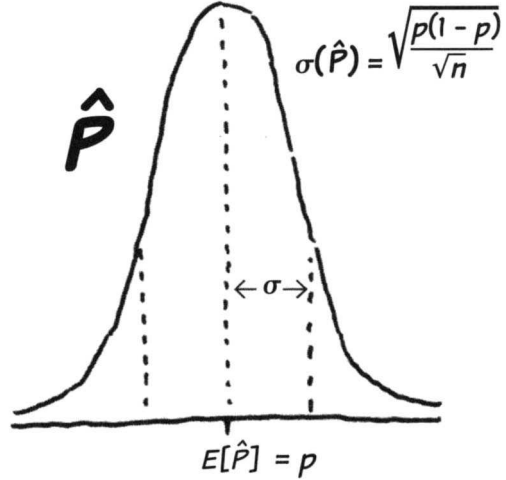
NOTE DE VOCABULAIRE : UNE **ESTIMATION** OU ESTIMATION PONCTUELLE EST UNE MESURE
OBSERVÉE SUR UN ÉCHANTILLON. UN **ESTIMATEUR** EST UNE RÈGLE POUR OBTENIR CES
ESTIMATIONS. L'ESTIMATEUR EST LA VARIABLE ALÉATOIRE $\hat{p} = X/n$.

LES STATISTIQUES IMPLIQUENT EN GÉNÉRAL UN PROCESSUS EN 4 ÉTAPES COMME NOUS VENONS DE LE FAIRE :

DÉFINIR LA POPULATION AVEC LE PARAMÈTRE INCONNU.



TROUVER UN ESTIMATEUR, SA DISTRIBUTION D'ÉCHANTILLONNAGE THÉORIQUE ET SON ERREUR-TYPE.



PRENDRE UN VRAI ÉCHANTILLON ALÉATOIRE ET TROUVER UNE ESTIMATION PONCTUELLE.



INDIQUER LE RÉSULTAT ET L'ERREUR STATISTIQUE D'ÉCHANTILLONNAGE.



Distribution d'échantillonnage de la MOYENNE

MAINTENANT, ON PASSE DES CLOUS EN LAITON AUX CORNICHONS...



LES FABRICANTS DE BOCAUX VEULENT CONNAÎTRE LA **TAILLE MOYENNE** DES CORNICHONS SANS AVOIR À INSPECTER TOUTES LES CUCURBITACÉES DE LA FRANCE. ILS SÉLECTIONNENT AU HASARD n CORNICHONS ET MESURENT LEURS TAILLES x_1, x_2, \dots, x_n .

À PRÉSENT, VOUS DEVEZ VOUS DOUTER QUE CHAQUE x_i EST UNE **VARIABLE ALÉATOIRE** EN TANT QUE RÉSULTAT D'UNE EXPÉRIENCE ALÉATOIRE.



SI μ EST LA MOYENNE (INCONNUE) ET σ EST L'ÉCART-TYPE DE LA **DISTRIBUTION DES TAILLES DE CORNICHONS**, ALORS :

$$E[x_i] = \mu$$

$$\sigma[x_i] = \sigma$$

POUR CHAQUE i (CAR x_i EST UNE TAILLE POSSIBLE DE CORNICHONS).



C'EST ÉTRANGE TOUT CE QU'ON SAIT SUR LES VARIABLES ALÉATOIRES, ALORS QU'IL Y A SEULEMENT UNE MINUTE ON NE SAVAIT PAS CE QU'ÉTAIT UNE VARIABLE ALÉATOIRE...

MAINTENANT REGARDONS LA MOYENNE D'ÉCHANTILLON, CELLE DES CORNICHONS SÉLECTIONNÉS. C'EST UNE NOUVELLE VARIABLE ALÉATOIRE DONNÉE PAR :

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

EST-CE QU'IL Y A QUELQUE CHOSE QUI N'EST PAS UNE VARIABLE ALÉATOIRE ?



COMME AUPARAVANT, NOUS VOUDRIONS CONNAÎTRE À QUEL POINT CETTE VALEUR APPROCHE μ . AUTREMENT DIT, SI ON DISPOSAIT DE PLUSIEURS ÉCHANTILLONS DIFFÉRENTS, QUELLE SERAIT LA DISTRIBUTION DE \bar{X} ? COMME NOUS CONNAISSONS $X_1, X_2 \dots$ ET X_n , ON SAIT AUSSI QUE :

$$E[\bar{X}] = \mu$$

$$\sigma(\bar{X}) = \frac{\sigma}{\sqrt{n}}$$

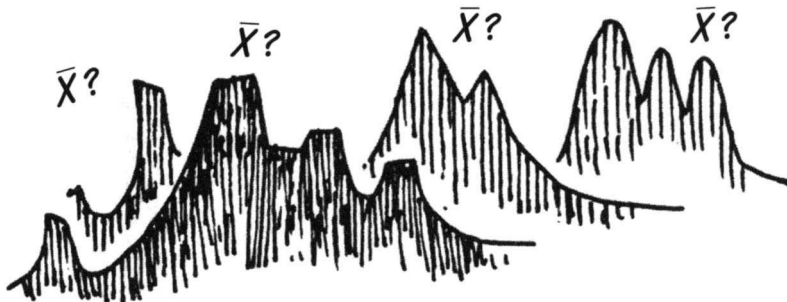
À NOUVEAU, ON RETROUVE NOTRE DÉNOMINATEUR MAGIQUE ! L'ÉCART-TYPE DES OBSERVATIONS D'ÉCHANTILLONS DÉPEND DE

$$\frac{1}{\sqrt{n}}$$



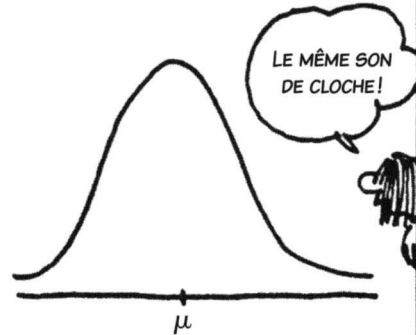
LES VARIANCES DES X_i/n S'ADDITIONNENT POUR DONNER LA VARIANCE DE \bar{X} .

MAIS NOUS NE CONNAISSONS PAS LA FORME DE LA DISTRIBUTION DE \bar{X} . LA **DISTRIBUTION D'ÉCHANTILLONNAGE** DE \hat{P} ÉTAIT PROCHE D'UNE NORMALE, CAR BASÉE SUR UNE VARIABLE ALÉATOIRE BINOMIALE. MAIS QU'EN EST-IL DE \bar{X} , L'ESTIMATEUR DE LA **MOYENNE** D'ÉCHANTILLON ?



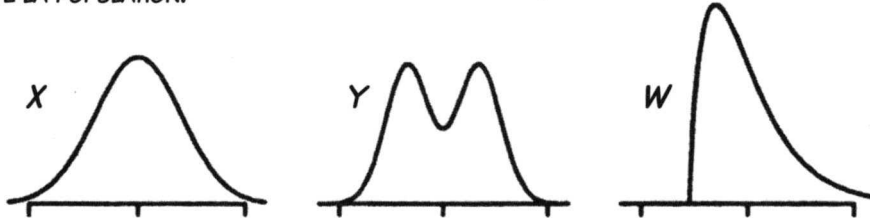
IL SE TROUVE QUE \bar{X} EST ÉGALEMENT PRESQUE NORMALE!
 CE RÉSULTAT BIEN CONNU PORTE AUSSI LE NOM
 DE **THÉORÈME CENTRAL LIMITE**.

IL DIT QUE SI L'ON PREND DES
 ÉCHANTILLONS ALÉATOIRES DE TAILLE
 n D'UNE POPULATION DE MOYENNE μ
 ET D'ÉCART-TYPE σ , ALORS PLUS n
 AUGMENTE, PLUS \bar{X} TEND VERS
 UNE **DISTRIBUTION NORMALE** DE
 MOYENNE μ ET D'ÉCART-TYPE σ/\sqrt{n} .
 AINSI,

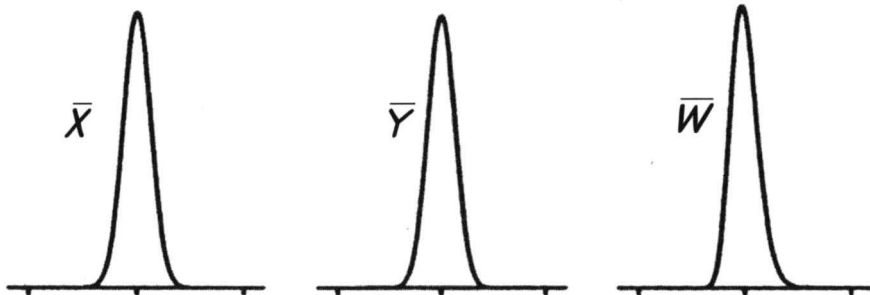


$$P(a \leq \bar{X} \leq b) \approx P\left(\frac{a - \mu}{\sigma/\sqrt{n}} \leq Z \leq \frac{b - \mu}{\sigma/\sqrt{n}}\right)$$

POURQUOI EST-CE REMARQUABLE? CELA INDIQUE QUE, QUELLE QUE SOIT LA FORME
 DE LA DISTRIBUTION INITIALE (DES TAILLES DE CORNICHONS POUR NOUS),
 LES **MOYENNES** D'ÉCHANTILLON TENDENT VERS UNE LOI **NORMALE**. POUR CONNAÎTRE
 LA DISTRIBUTION DE \bar{X} , IL NOUS SUFFIT DE CONNAÎTRE LA MOYENNE ET L'ÉCART-TYPE
 DE LA POPULATION.



LES TROIS DENSITÉS PRÉCÉDENTES ONT LES MÊMES MOYENNES ET ÉCARTS-TYPES. MALGRÉ
 LEURS FORMES DIFFÉRENTES, LORSQUE $n = 10$, LES DISTRIBUTIONS D'ÉCHANTILLONNAGE
 DE LA MOYENNE \bar{X} SONT QUASIMENT IDENTIQUES.



Le «t de Student»

AUSSI INCROYABLE QUE SOIT LE THÉORÈME CENTRAL LIMITE, IL A AU MOINS DEUX DÉFAUTS.



LE PREMIER : IL FAUT QUE n SOIT GRAND.

LE SECOND : IL FAUT CONNAÎTRE L'ÉCART-TYPE σ .

MAIS LES ÉCHANTILLONS SONT SOUVENT PETITS ET σ EST GÉNÉRALEMENT INCONNU. QUOI QU'IL EN SOIT DANS LE CAS DES CORNICHONS, NOUS N'AVONS **AUCUNE IDÉE** DE LA DISPERSION DES TAILLES DES CORNICHONS AUTOUR DE LA MOYENNE.



ON PEUT ALORS **ESTIMER** σ EN PRENANT **s L'ÉCART-TYPE À L'INTÉRIEUR DE L'ÉCHANTILLON** OBTENU À PARTIR DE :

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

AINSI À L'INTÉRIEUR DE LA VARIABLE ALÉATOIRE :

$$z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

ON REMPLACE σ PAR s ET ON DÉFINIT UNE NOUVELLE **VARIABLE ALÉATOIRE t** :

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}}$$



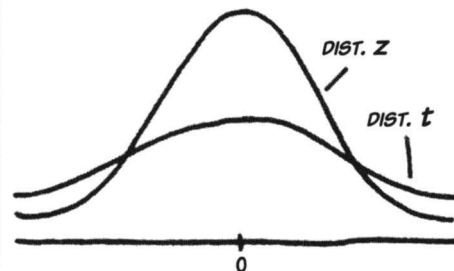
ON PEUT VOIR LA VARIABLE ALÉATOIRE t COMME LA MEILLEURE APPROXIMATION POSSIBLE ÉTANT DONNÉ LES CIRCONSTANCES. SA DISTRIBUTION EST APPELÉE t DE STUDENT, CAR ELLE A ÉTÉ PUBLIÉE PAR SON INVENTEUR **WILLIAM GOSSET** SOUS LE PSEUDONYME «ÉTUDIANT».



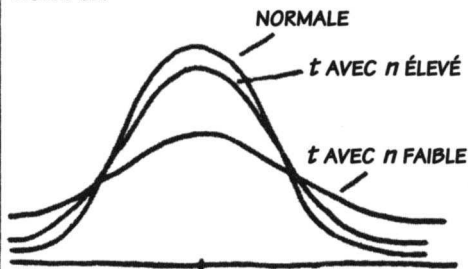
EN FAISANT L'HYPOTHÈSE QUE LA **POPULATION ORIGINALE DES OBSERVATIONS ÉTAIT NORMALE**, OU PRESQUE, «STUDENT» PUT FOURNIR UNE CONCLUSION :



LE t DE STUDENT EST PLUS DISPERSÉ ET PLUS APLATI QU'UNE LOI NORMALE (z). CELA S'EXPLIQUE PAR L'UTILISATION DE L'ÉCART-TYPE D'ÉCHANTILLON (s) QUI INTRODUIT ENCORE PLUS D'INCERTITUDE DANS LA DISTRIBUTION.



LA QUANTITÉ DE DISPERSION DÉPEND DE n , LA TAILLE DE L'ÉCHANTILLON. PLUS LA TAILLE D'ÉCHANTILLON EST ÉLEVÉE, PLUS s APPROXIME BIEN σ ET t SE RAPPROCHE DE z , LA LOI NORMALE.



GOSSET PUT ÉTABLIR DES TABLES DE t POUR DIFFÉRENTES TAILLES D'ÉCHANTILLON ET NOUS LES UTILISERONS DANS LE PROCHAIN CHAPITRE.



DANS CE CHAPITRE, NOUS AVONS VU LE PROBLÈME CENTRAL DES **STATISTIQUES APPLIQUÉES** : COMMENT CONSTRUIRE UN **ÉCHANTILLON** D'UNE GRANDE POPULATION AFIN QUE LES ANALYSES STATISTIQUES SOIENT VALIDES. OUTRE LA « NORME STANDARD » DE L'ÉCHANTILLON ALÉATOIRE SIMPLE, NOUS AVONS VU D'AUTRES MÉTHODES D'ÉCHANTILLONNAGE QUI PEUVENT ÊTRE UTILISÉES POUR DES RAISONS PRATIQUES D'EFFICACITÉ OU DE COÛT.



ENSUITE, ET AVEC UN ÉCHANTILLON ALÉATOIRE SIMPLE, NOUS AVONS VU COMMENT LES STATISTIQUES D'ÉCHANTILLONS ÉTAIENT **DISTRIBUÉES**. AINSI ON A VU QUE PRENDRE UN ÉCHANTILLON CONSTITUAIT UNE **EXPÉRIENCE ALÉATOIRE**, ET QUE LES STATISTIQUES OBSERVÉES DEVENAIENT DES **VARIABLES ALÉATOIRES**.



ON A VU QUE LA VARIABLE ALÉATOIRE **PROPORTION D'ÉCHANTILLON \hat{p}** , ET QUE LA VARIABLE ALÉATOIRE **MOYENNE D'ÉCHANTILLON \bar{X}** ÉTAIENT À PEU PRÈS DISTRIBUÉES NORMALEMENT LORSQUE LA TAILLE D'ÉCHANTILLON ÉTAIT SUFFISANTE. ON A AUSSI INTRODUIT LES VARIABLES **STANDARDS z ET t** , QUI SERVIRONT SUR LES RÉSULTATS D'ÉCHANTILLONS POUR LES CHAPITRES SUIVANTS.



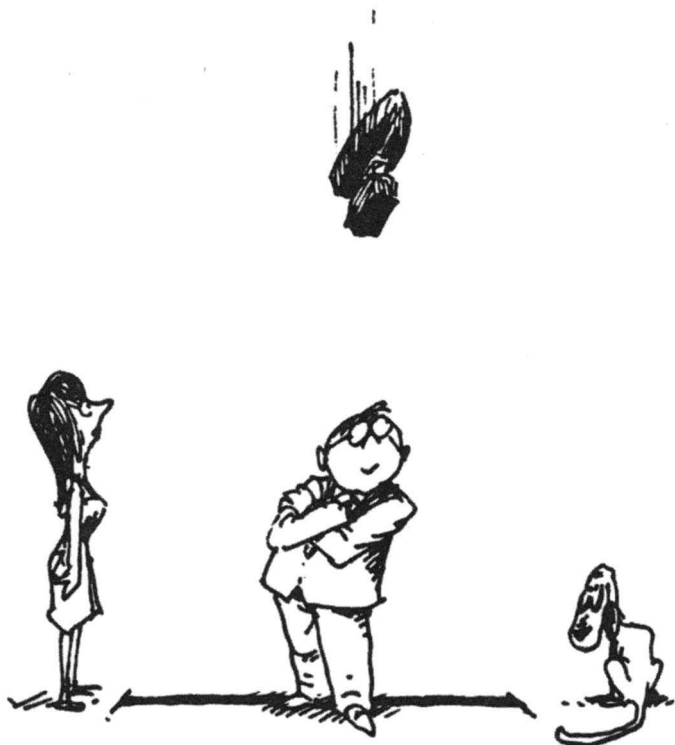
*AVEC DU THÉ BIEN SÛR!

DANS LES DEUX PROCHAINS CHAPITRES, NOUS ALLONS
VOIR COMMENT UTILISER CES DISTRIBUTIONS POUR FAIRE
DES **INFÉRENCES STATISTIQUES** : À PARTIR D'UN SEUL
ÉCHANTILLON, UN SONDAGE D'OPINION PAR EXEMPLE, COMMENT
MOBILISER NOS CONNAISSANCES SUR \hat{P} ET \bar{X} POUR LES ÉVALUER ?

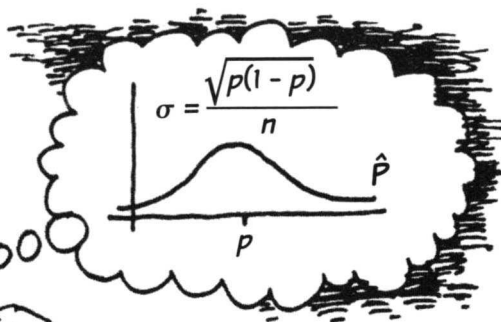


Chapitre 7

Intervalles de confiance



DANS LE CHAPITRE PRÉCÉDENT,
NOUS AVONS VU LES MÉTHODES
D'ÉCHANTILLONNAGE. À PARTIR
D'UNE GRANDE POPULATION,
NOUS AVONS ANALYSÉ COMMENT
LA DISTRIBUTION D'ESTIMATEURS
D'ÉCHANTILLON VARIAIT SELON
LES ÉCHANTILLONS.



DANS CE CHAPITRE, NOUS FERONS L'INVERSE. À PARTIR D'UN ÉCHANTILLON,
NOUS NOUS DEMANDERONS QUEL EST LE SYSTÈME ALÉATOIRE QUI A GÉNÉRÉ
NOS STATISTIQUES ?



AINSI,
À PARTIR D'UNE SEULE BOÎTE
DE CLOUS ET ÉTANT DONNÉ
LES RÉSULTATS
DU CHAPITRE PRÉCÉDENT,
QUE POUVONS-NOUS
CONCLURE ?

CELA CHANGE NOTRE MODE DE PENSÉE :
ON PASSE D'UN RAISONNEMENT DÉDUCTIF
À UN RAISONNEMENT **INDUCTIF**.



C'EST COMME UNE ENQUÊTE
CRIMINELLE, WATSON!



DANS UN RAISONNEMENT **DÉDUCTIF**,
ON PART D'UNE HYPOTHÈSE POUR ARRIVER
À UNE CONCLUSION : « SI LORD FINE-GACHETTE
A COMMIS LE CRIME, ALORS IL A EFFACÉ
SES EMPREINTES DE L'ARME. »

LE RAISONNEMENT **INDUCTIF**
EST À **REBOURS**, ON PART
D'UN ENSEMBLE D'OBSERVATIONS
VERS UNE HYPOTHÈSE VRAISEMBLABLE.

HUM!
LE **MONOGRAMME**
DE LORD FINE-GACHETTE
EST SUR CE **MOUCHOIR**
ET CE **REVOLVER**. FINE-GACHETTE
EST LE MEURTRIER, WATSON,
J'EN SUIS **SÛR À 95 %**!



INDUCTION
BRILLANTE,
HOLMES!



À BIEN DES ÉGARDS, LA SCIENCE ET DONC LES STATISTIQUES RESSEMBLENT
À UN TRAVAIL DE DÉTECTIVE. À PARTIR D'UN ENSEMBLE D'OBSERVATIONS,
NOUS NOUS INTERROGEONS SUR CE QUI LES A GÉNÉRÉES.

ESTIMATION d'intervalles de confiance

C'EST UNE FORME D'INFÉRENCE STATISTIQUE,
QUE L'ON VOIT APPARAÎTRE CHAQUE FOIS
PENDANT LES PÉRIODES ÉLECTORALES...



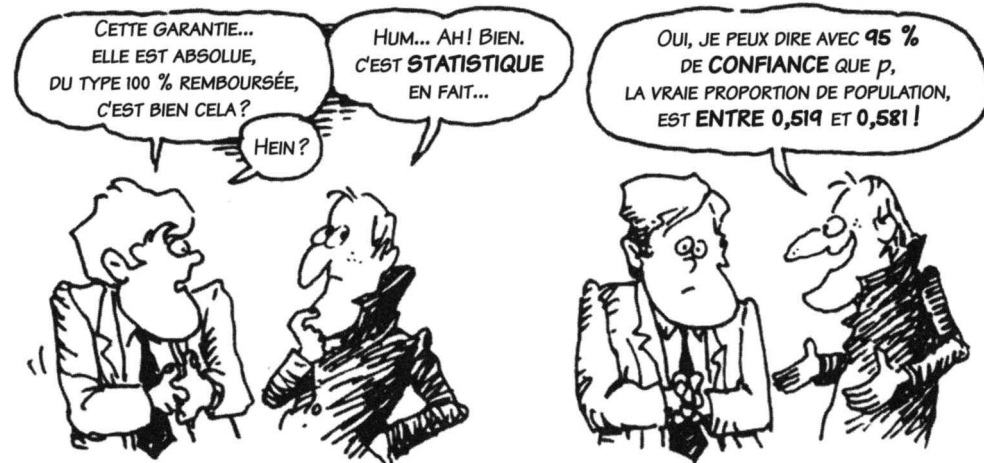
QUELQUE PART, LORS D'UNE ÉLECTION RÉCENTE, LE SÉNATEUR EN PLACE **ASTUTE** COMMANDE UN SONDAGE À **HOLMES RECHERCHE INSTITUT**. LE SONDEUR HOLMES PREND ALORS UN ÉCHANTILLON ALÉATOIRE SIMPLE DE 1000 ÉLECTEURS ET LEUR DEMANDE CE QU'ILS PENSENT D'ASTUTE.



APRÈS AVOIR CENSURÉ LES REMARQUES DES QUELQUES VALEURS EXTRÊMEMENT GRINCHEUSES, HOLMES TROUVE QUE **550** VOTANTS SONT EN FAVEUR DE SON CLIENT, LE SÉNATEUR **ASTUTE**.



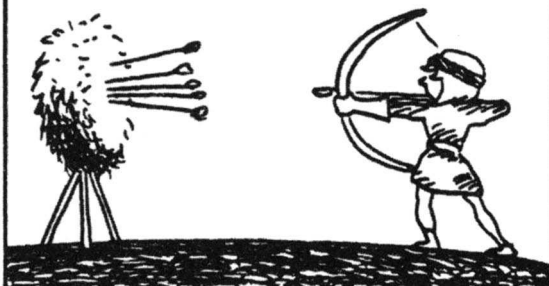
CELA CONSTITUE
L'ESTIMATION PONCTUELLE.



LE SÉNATEUR ASTUTE DEMEURE PERPLEXE! ALORS HOLMES LUI DONNE UNE **leçon de tir à l'arc.**



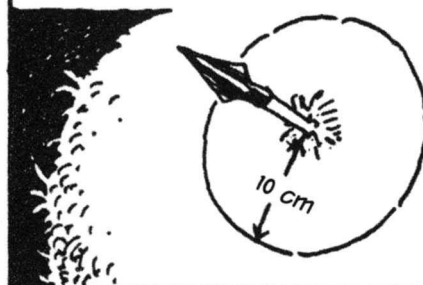
CONSIDÉRONS UN ARCHER SONDEUR QUI VISE UNE CIBLE. SUPPOSONS QU'IL TOUCHE 95 % DU TEMPS LE CENTRE DE LA CIBLE DANS UN RAYON DE 10 cm. AINSI, SEULE UNE FLÈCHE SUR 20 N'ATTEINT PAS CE DISQUE.



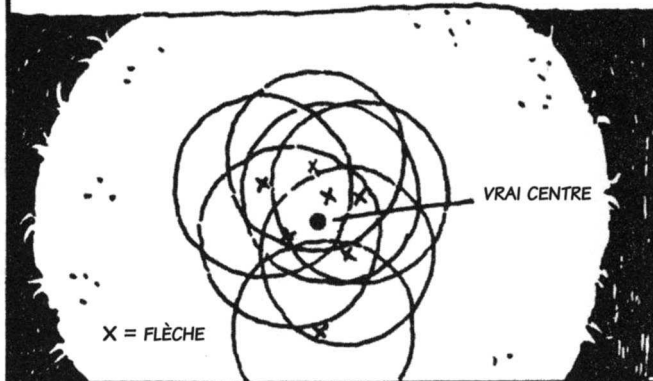
NOTRE BRAVE DÉTECTIVE EST ASSIS DERRIÈRE LA CIBLE ET N'EN VOIT PAS LE CENTRE. L'ARCHER TIRE UNE SEULE FLÈCHE.



CONNAISSANT L'AGILITÉ DU TIREUR, LE DÉTECTIVE TRACE UN CERCLE DE RAYON 10 cm AUTOUR DE LA POINTE DE LA FLÈCHE. IL EST MAINTENANT **CONFIANT À 95 %** QUE LE DISQUE INCLUT LE CENTRE DE LA CIBLE.



SON RAISONNEMENT EST QUE, S'IL TRACE DES CERCLES DE RAYON DE 10 cm AUTOUR DE **PLUSIEURS FLÈCHES**, LES DISQUES VONT INCLURE LE CENTRE DANS 95 % DES CAS.



(LES PROBABILISTES UTILISENT LE TERME **STOCHASTIQUE** POUR DÉCRIRE CES MODÈLES ALÉATOIRES. CELA VIENT DU GREC **STOKHAZETHAI** QUI VEUT DIRE VISER UNE CIBLE, ET DE **STOKHOS** UNE FLÈCHE.)



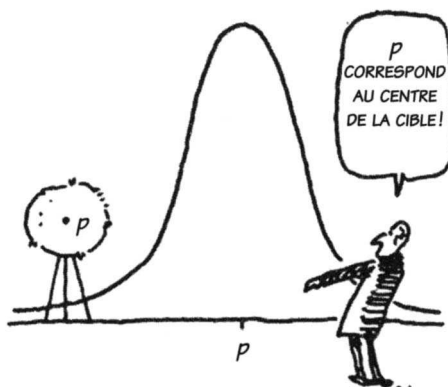


HOLMES TRADUIT MAINTENANT LA LEÇON DE TIR À L'ARC DANS LE LANGAGE DÉVELOPPÉ AU CHAPITRE PRÉCÉDENT.

Premier temps. TIREZ BEAUCOUP DE FLÈCHES.

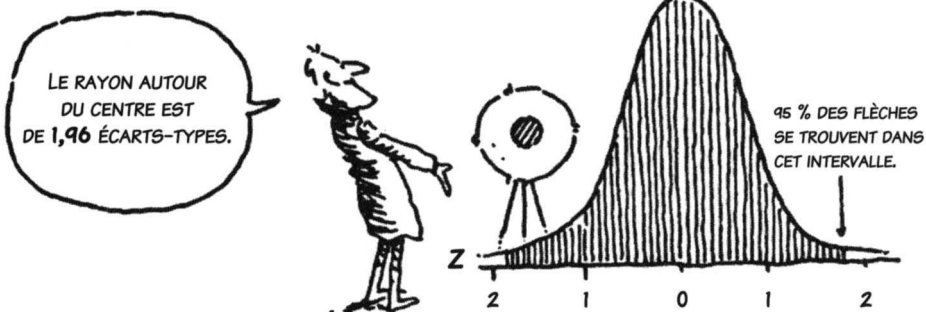
UN CALCUL DE PROBABILITÉ PERMET D'ESTIMER LA POSITION DU CENTRE DE LA CIBLE. LES ESTIMATIONS PONCTUELLES \hat{p} CORRESPONDENT À NOS FLÈCHES. LA DISTRIBUTION D'ÉCHANTILLONNAGE DE \hat{p} EST PRESQUE NORMALE DE MOYENNE μ ET D'ÉCART-TYPE

$$\sigma(\hat{p}) = \frac{\sqrt{p(1-p)}}{\sqrt{n}}$$



COMME LA COURBE EST NORMALE, NOUS UTILISONS SON Z-SCORE ET UNE TABLE STANDARD POUR ÉVALUER LA LARGEUR DE L'INTERVALLE DANS LEQUEL 95 % DES FLÈCHES TOMBENT (NOUS ALLONS VOIR EXACTEMENT COMMENT CALCULER CELA DANS QUELQUES PAGES). NOUS TROUVONS QUE CETTE LARGEUR EST DE 1,96 ÉCARTS-TYPES.

$$0,95 = P(-1,96 \leq z \leq 1,96)$$



MAINTENANT, FAISONS UN PEU D'ALGÈBRE.
ON TRANSFORME EN LOI NORMALE STANDARD Z :

$$0,95 \approx P\left(-1,96 \leq \frac{\hat{P} - p}{\sigma(\hat{P})} \leq 1,96\right)$$

CE QUI SE TRADUIT PAR :

$$0,95 \approx P(p - 1,96\sigma(\hat{P}) \leq \hat{P} \leq p + 1,96\sigma(\hat{P}))$$

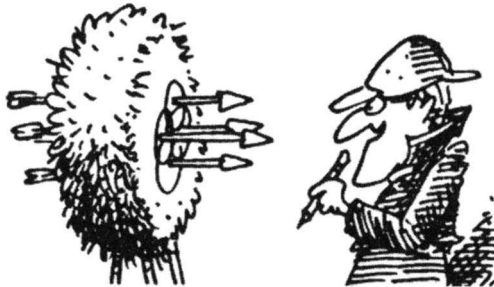
CE QUI REVIENT À DIRE QUE 95 % DES « FLÈCHES » \hat{p} TOMBENT ENTRE $p - 1,96\sigma(\hat{P})$ ET $p + 1,96\sigma(\hat{P})$.



MAINTENANT, NOUS POUVONS REGARDER LA CIBLE DE DERRIÈRE. ENCORE UN TOUR DE MANIVELLE ALGÈBRIQUE ET ON OBTIENT :

$$0,95 \approx P(\hat{P} - 1,96\sigma(\hat{P}) \leq p \leq \hat{P} + 1,96\sigma(\hat{P}))$$

ICI, ON DESSINE LES CERCLES AUTOUR DES POINTES DE FLÈCHE, AINSI ON TRACE DES INTERVALLES AUTOUR DE \hat{p} , ET ON DIT QUE 95 % D'ENTRE EUX COUVRENT p .



MAIS IL Y A UN LÉGER PROBLÈME : ON NE CONNAÎT PAS EXACTEMENT LA TAILLE AUTOUR DE LA CIBLE, CAR SA LARGEUR EST UN MULTIPLE DE $\sigma(\hat{P})$ QUI DÉPEND DE L'INCONNU p .



MAINTENANT, LES DISQUES ONT DES RAYONS DIFFÉRENTS MAIS CE N'EST PAS GRAVE, VRAIMENT...

ON ESQUIVE DONC LE PROBLÈME EN UTILISANT L'ESTIMATION PONCTUELLE \hat{p} À LA PLACE DE p POUR DÉTERMINER L'ERREUR-TYPE ESTIMÉE OU ERREUR-STANDARD :

$$s(\hat{P}) = \frac{\sqrt{\hat{p}(1 - \hat{p})}}{\sqrt{n}}$$

C'EST SUFFISAMMENT PROCHE...
C'EST AU MIEUX... ET CELA PEUT SE JUSTIFIER THÉORIQUEMENT !

MAINTENANT LA FORMULE DEVIENT :

$$0,95 \approx P(\hat{p} - 1,96s(\hat{p}) \leq p \leq \hat{p} + 1,96s(\hat{p}))$$

DE NOUVEAU, CETTE ÉQUATION DÉCRIT LA PROBABILITÉ QUE LA VRAIE VALEUR DE PROPORTION DE LA POPULATION SOIT DANS L'INTERVALLE ALÉATOIRE.

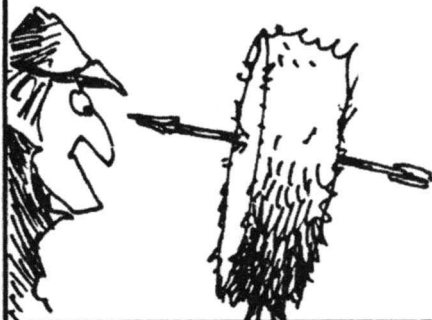
$$[\hat{p} - 1,96s(\hat{p}), \hat{p} + 1,96s(\hat{p})]$$

SI L'ON PREND DES ÉCHANTILLONS RÉPÉTÉS, LES INTERVALLES $[\hat{p} - 1,96s(\hat{p}), \hat{p} + 1,96s(\hat{p})]$ COUVRIRONT p DANS 95 % DES CAS.



MAINTENANT QUE LES CALCULS SONT FAITS, IL EST TEMPS DE PASSER AU...

Second temps. LE TRAVAIL DE DÉTECTIVE. DANS UN VRAI SONDAGE, HOLMES PREND UN SEUL ÉCHANTILLON DE 1000 VOTANTS, TROUVE $\hat{p} = 0,55$ ET VEUT EN DÉDUIRE p .



IL UTILISE NOTRE PREMIER TEMPS POUR CALCULER :

$$s(\hat{p}) = \frac{\sqrt{(0,55) \times (0,45)}}{\sqrt{1000}} = 0,0157$$

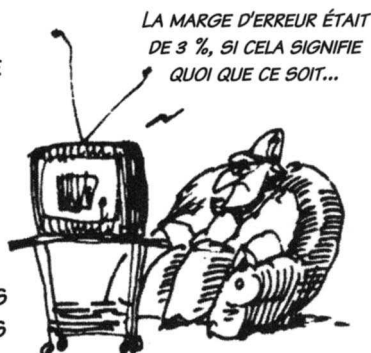
ET CONCLUT QU'AVEC 95 % DE CONFIANCE p EST DANS L'INTERVALLE :

$$\begin{aligned} & \hat{p} \pm 1,96s(\hat{p}) \\ &= 0,550 \pm 1,96 \times 0,0157 \\ &= 0,550 \pm 0,031 \end{aligned}$$

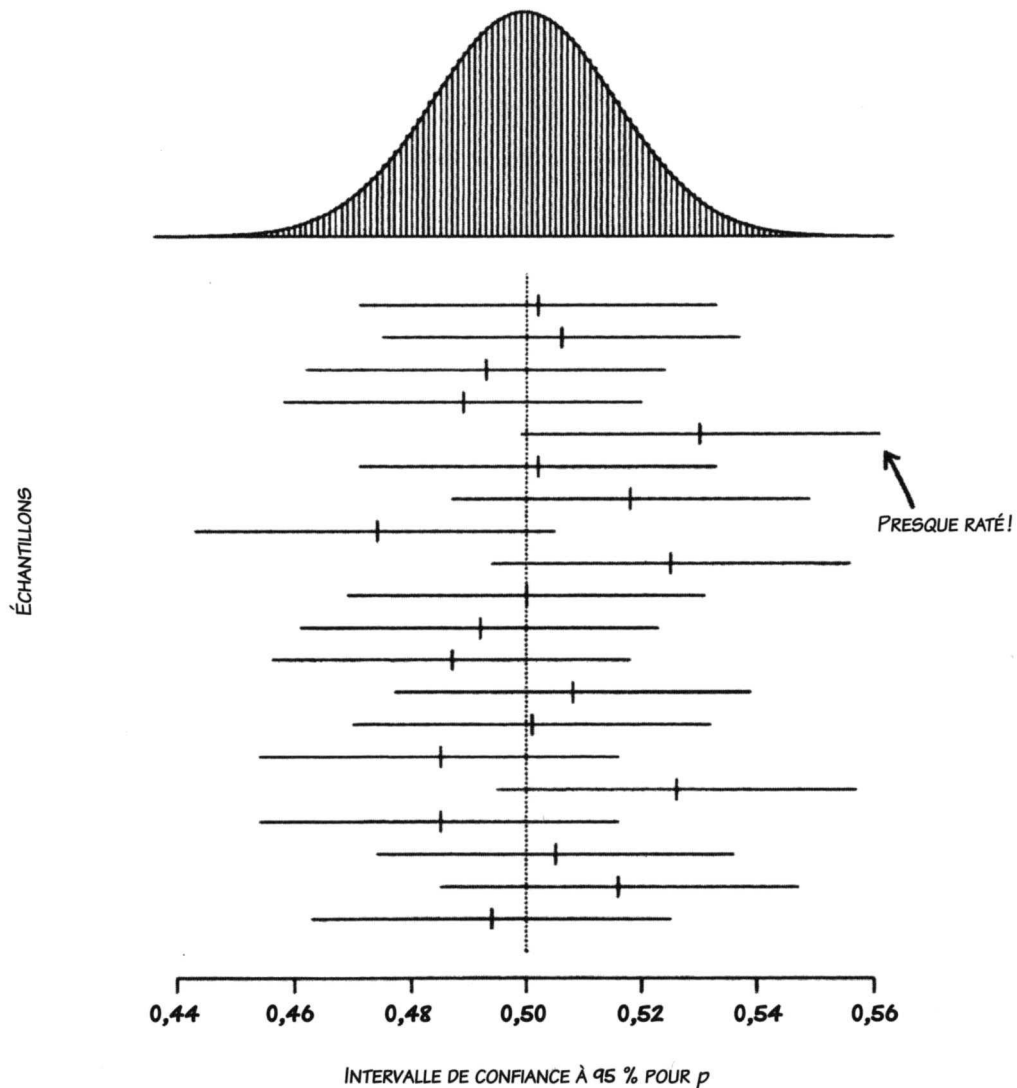
C'EST CE À QUOI FONT RÉFÉRENCE LES SONDAGES EN PARLANT DE « MARGE D'ERREUR ». DANS NOTRE CAS, HOLMES A TROUVÉ QUE :

$$0,519 \leq p \leq 0,581$$

EN D'AUTRES TERMES, $p = 55\%$ AVEC 3 % DE MARGE D'ERREUR (LES SONDAGES UTILISENT TYPIQUEMENT DES INTERVALLES À 95 %).



NOUS MONTRONS SUR CETTE PAGE LES RÉSULTATS D'UNE SIMULATION PAR ORDINATEUR DE VINGT ÉCHANTILLONS DE TAILLE $n = 1000$. LA SIMULATION UTILISE LA VALEUR $p = 0,5$. LA DISTRIBUTION D'ÉCHANTILLONNAGE DE \hat{P} (NORMALE DE MOYENNE p ET ÉCART-TYPE $\sigma = \sqrt{p(1-p)/n}$) APPARAÎT EN HAUT. EN DESSOUS, IL Y A LES INTERVALLES DE CONFIANCE À 95 % DE CHAQUE ÉCHANTILLON. UN SUR VINGT (SOIT 5 %) DE CES INTERVALLES **NE COUVRIRONT PAS LA VALEUR** $p = 0,5$.

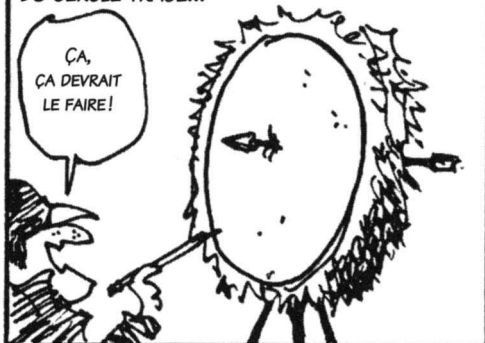


BIEN QUE 95 %
DE CONFIANCE SUFFISENT
POUR LES SONDAGES
DES JOURNAUX,
CE N'EST PAS ASSEZ
POUR LE SÉNATEUR
ASTUTE. IL VEUT 99 %!

SI C'EST INFÉRIEUR,
MES SPONSORS FORTUNÉS
N'INVESTIRONT PAS – JE VEUX DIRE
NE **CONTRIBUERONT PAS** –
À MA LUTTE POUR
LA LIBERTÉ
ET LA JUSTICE.



COMMENT AUGMENTER LA CONFIANCE?
PAR ANALOGIE AVEC LE TIR À L'ARC, DEUX
MÉTHODES SONT POSSIBLES : LA PREMIÈRE
CONSISTE À AUGMENTER LA TAILLE
DU CERCLE TRACÉ...



ET UNE AUTRE FAÇON SERAIT D'AMÉLIORER
LA VISÉE INITIALE DE L'ARCHER, DE MANIÈRE
QUE LE TIR DES FLÈCHES SOIT PLUS GROUPE
AU CENTRE DE LA MIRE.



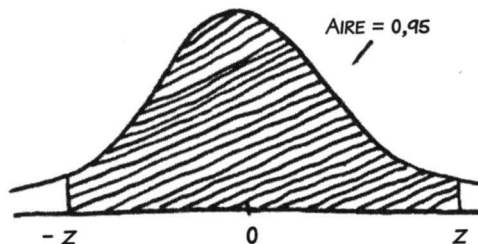
LA PREMIÈRE MÉTHODE REVIENT À ÉLARGIR L'INTERVALLE DE CONFIANCE. PLUS LA MARGE
D'ERREUR EST IMPORTANTE, PLUS ON EST SÛR QUE LA VRAIE VALEUR p SERA DANS L'INTERVALLE.



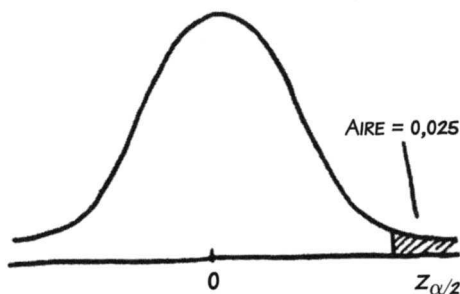
PEUT-ÊTRE EST-IL TEMPS DE VOIR EXACTEMENT
COMMENT ON DÉTERMINE LES BORNES
DE CES INTERVALLES DE CONFIANCE...

ICI, LE NOMBRE PERTINENT S'APPELLE α (ALPHA). IL MESURE LA DIFFÉRENCE ENTRE LE SEUIL DE CONFIANCE VOULU ET LA CERTITUDE. PAR EXEMPLE, QUAND LE SEUIL DE CONFIANCE VAUT 95 %, OU 0,95, ALORS $\alpha = 0,05$. ON PARLE DONC D'UN SEUIL (OU NIVEAU) DE CONFIANCE DE $(1 - \alpha)$.

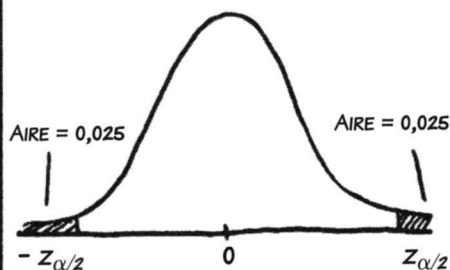
CHERCHER UN INTERVALLE DE CONFIANCE À $(1 - \alpha)$ REVIENT À REGARDER LA LOI NORMALE ET À CHERCHER LES POINTS $\pm z$ ENTRE LESQUELS L'AIRE VAUT $(1 - \alpha)$.



LE POINT $z_{\alpha/2}$, APPELÉ **FRACTILE**, EST LA VALEUR z AU-DESSUS DE LAQUELLE L'AIRE EST DE $\alpha/2 = 0,025$.



ON ENLÈVE LES QUEUES DE DISTRIBUTION AUX EXTRÉMITÉS DE LA COURBE, QUI ONT UNE AIRE TOTALE DE $\alpha = \alpha/2 + \alpha/2$.



Z	-2,5	-2,4	-2,3	-2,2	-2,1	-2,0	-1,9	-1,8	-1,7	-1,6	-1,5	-1,4
F(z)	0,006	0,008	0,011	0,014	0,018	0,023	0,029	0,036	0,045	0,055	0,067	0,081

ON PEUT TROUVER $z_{\alpha/2}$ DIRECTEMENT AVEC LA TABLE DE LOI NORMALE DE LA PAGE 84. IL S'AGIT DU POINT AYANT CETTE PROPRIÉTÉ :

$$P(z \geq z_{\alpha/2}) = \frac{\alpha}{2}$$

EN PARTICULIER :

$$P(z \geq z_{0,025}) = 0,025$$



VOICI UNE PETITE TABLE DE FRACTILES
POUR DIFFÉRENTS SEUILS DE CONFIANCE :

$1 - \alpha$	0,8	0,9	0,95	0,99
α	0,2	0,1	0,05	0,01
$\alpha/2$	0,1	0,05	0,025	0,005
$z_{\alpha/2}$	1,28	1,64	1,96	2,58



POUR FAIRE UN INTERVALLE DE CONFIANCE À 99 %, ON UTILISE CETTE TABLE POUR ÉCRIRE :

$$0,99 = P(\hat{p} - 2,58s(\hat{p}) \leq p \leq \hat{p} + 2,58s(\hat{p}))$$

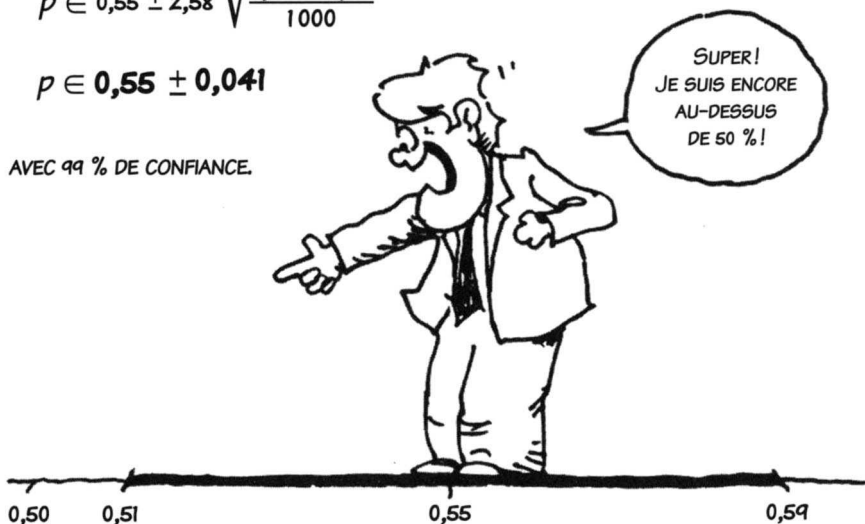
QUE NOUS ABRÉGEONS :

$$p \in \hat{p} \pm 2,58s(\hat{p})$$

$$p \in 0,55 \pm 2,58 \sqrt{\frac{0,55 \times 0,45}{1000}}$$

$$p \in \mathbf{0,55 \pm 0,041}$$

AVEC 99 % DE CONFIANCE.



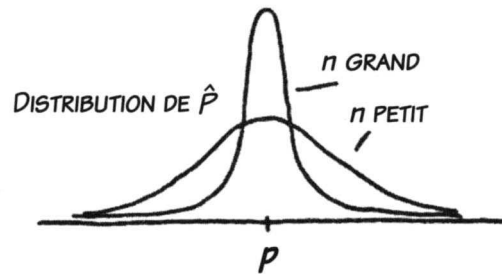
AUGMENTER LES INTERVALLES EST UNE FAÇON D'AUGMENTER LA CONFIANCE DANS SES RÉSULTATS. COMME NOUS L'AVONS DIT, UNE AUTRE FAÇON SERAIT DE VISER PLUS JUSTE AVEC NOS FLÈCHES. SI NOUS SAVIONS QUE 95 % DE NOS FLÈCHES SONT À 1 CM DE LA CIBLE, NOS ESTIMATIONS POURRAIENT ÊTRE PLUS AFFÛTÉES.



COMMENT POUVONS-NOUS FAIRE? EN AUGMENTANT LA TAILLE DE L'ÉCHANTILLON! LA LARGEUR DE L'INTERVALLE DE CONFIANCE DÉPEND DE CETTE TAILLE : L'INTERVALLE EST DE LA FORME $\hat{p} \pm E$ OÙ L'ERREUR E EST DONNÉE PAR :

$$E = z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

PLUS n EST GRAND, PLUS L'ERREUR EST PETITE (PAR EXEMPLE QUADRUPLER n RÉDUIT DE MOITIÉ LA LARGEUR DE L'INTERVALLE).



ASTUTE DEMANDE À HOLMES UNE PETITE ERREUR AVEC UNE GRANDE CONFIANCE – DISONS 99 % DE CONFIANCE AVEC $E = \pm 0,01$. HOLMES CALCULE DONC n :

$$n = \frac{(z_{\alpha/2})^2 p^*(1 - p^*)}{E^2}$$

ICI p^* EST UNE PRÉVISION DE LA VRAIE PROPORTION p (RAPPELEZ-VOUS QUE CELA SE FAIT AVANT LA MESURE D'ÉCHANTILLON).



EN PRENANT UNE PRÉVISION PRUDENTE
DE $p^* = 0,5$, HOLMES TROUVE :

$$n = \frac{(2,58)^2 (0,5)^2}{(0,01)^2}$$
$$= \frac{(6,65)(0,25)}{(0,0001)}$$
$$= 16\,641$$

1 000 VOTANTS DONNAIENT UNE ERREUR
DE 3 % AVEC 95 % DE CONFIANCE.
POUR AVOIR 1 % D'ERREUR AVEC 99 %
DE CONFIANCE, HOLMES DOIT SONDER
16 641 VOTANTS.



ILS FONT ALORS LE SONDAGE
ET VONT À L'ÉLECTION AVEC
99 % DE CONFIANCE.

MAIS... TOUTES CES PROBABILITÉS SONT UTILES AVANT LES ÉLECTIONS.
APRÈS L'ÉLECTION, LE SÉNATEUR EST SOIT 100 % ÉLU, SOIT 100 % PERDANT...
ET MALGRÉ TOUT, LE SÉNATEUR ASTUTE **PERD** L'ÉLECTION.



QUE S'EST-IL
PASSÉ ?

IL S'EST PASSÉ QUE LES POLITICIENS NE SONT PAS ÉLUS AVEC LES SONDAGES !



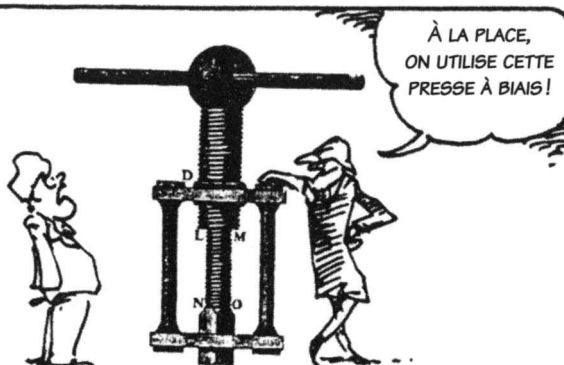
QUELQUES PROBLÈMES DES SONDAGES PAR RAPPORT AUX ÉLECTIONS :



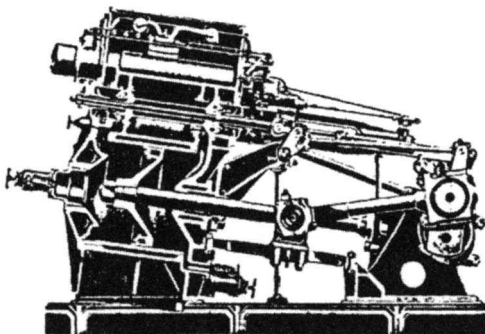
IL N'Y A AUCUN MOYEN POUR UN SONDEUR D'ENTRER DANS LA TÊTE D'UN VOTANT, DE SAVOIR S'IL VA VOTER, S'IL MENT OU S'IL VA CHANGER D'AVIS AVANT L'ÉLECTION. DE GRANDS ÉCHANTILLONS NE PEUVENT PAS RÉDUIRE CE TYPE D'ERREUR.



COMME CES ERREURS PEUVENT ÊTRE IMPORTANTES, ON GAGNE RAREMENT À PRENDRE DE TRÈS GRANDS ÉCHANTILLONS.



DANS CINQ ÉLECTIONS PRÉSIDENTIELLES CONSÉCUTIVES, LES SONDAGES GALLUP ONT INTERVIEWÉ MOINS DE 4 000 VOTANTS POUR CHAQUE ÉLECTION. POURTANT DANS CES CINQ ÉLECTIONS, L'ERREUR DE PRÉDICTION DU RÉSULTAT CALCULÉ PAR GALLUP ÉTAIT DE MOINS DE 2 %.



CE SUCCÈS EST DÙ À L'UTILISATION D'ESTIMATEURS QUI PRENNENT EN COMPTE LES NON-RÉPONSES ET FILTRENT LES ÉLECTEURS ÉLIGIBLES QUI N'IRONT VRAISEMBLABLEMENT PAS VOTER.



EN RÉSUMÉ, PROPORTION ESTIMÉE
= VRAIE PROPORTION + BIAIS + ERREUR
ALÉATOIRE D'ÉCHANTILLON. MÊME LES SONDEURS
ONT DES RESSOURCES LIMITÉES. ILS CHOISSENT
JUDICIEUSEMENT DE DÉPENSER DE L'ARGENT POUR
RÉDUIRE LES BIAIS PLUTÔT QUE D'AUGMENTER
LE NOMBRE DE VOTANTS SONDÉS À PLUS DE 4 000.

Intervalle de confiance pour μ

JUSQU'À PRÉSENT, NOUS AVONS EXAMINÉ LES INTERVALLES DE CONFIANCE POUR UNE PROPORTION p D'UNE POPULATION. LE MÊME TYPE DE RAISONNEMENT FONCTIONNE POUR UNE MOYENNE μ DE POPULATION.



DANS LE CHAPITRE PRÉCÉDENT, NOUS AVONS VU PAGE 105 QUE LA DISTRIBUTION D'ÉCHANTILLONNAGE DES MOYENNES \bar{X} EST PRESQUE NORMALE, CENTRÉE SUR LA MOYENNE DE POPULATION μ ET D'ÉCART-TYPE σ/\sqrt{n} OÙ σ EST L'ÉCART-TYPE DE LA POPULATION. AINSI, SI n EST GRAND :

$$0,95 = P(-1,96 \leq Z \leq 1,96) \\ \approx P(-1,96 \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq 1,96)$$

À NOUVEAU, NE CONNAISSANT PAS σ , ON REMPLACE σ PAR s , L'ÉCART-TYPE ÉCHANTILLON, ET ON OBTIENT :

$$0,95 \approx P(-1,96 \leq \frac{\bar{X} - \mu}{s/\sqrt{n}} \leq 1,96)$$



LE TERME s/\sqrt{n} EST APPELÉ L'ERREUR-TYPE ESTIMÉE ET SE NOTE $s(\bar{X})$. ON EN CONCLUT QUE :

$$0,95 \approx P(\bar{X} - 1,96s(\bar{X}) \leq \mu \leq \bar{X} + 1,96s(\bar{X}))$$

OÙ :

$$s(\bar{X}) = \frac{s}{\sqrt{n}}$$



COMME PRÉCÉDEMMENT, NOUS AVONS
TROUVÉ QUE L'INTERVALLE ALÉATOIRE

$$\bar{X} \pm 1,96s(\bar{X})$$

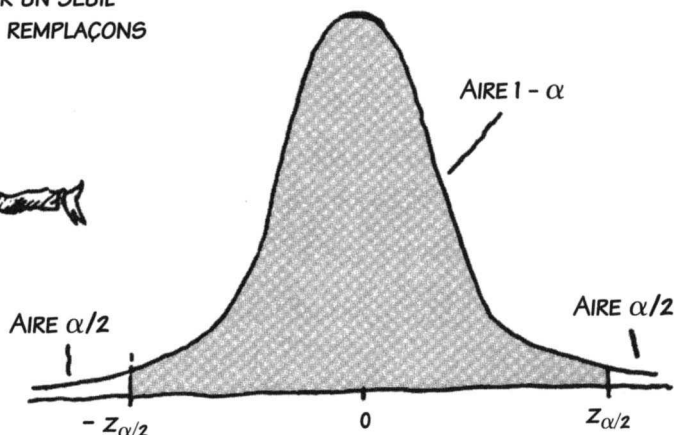
RECOUVRE LA VRAIE MOYENNE μ
AVEC UNE PROBABILITÉ DE 0,95...
DONC MAINTENANT ON PEUT FAIRE APPEL
À SHERLOCK HOLMES POUR EFFECTUER
L'**INFÉRENCE STATISTIQUE** FONDÉE
SUR UN **SEUL** ÉCHANTILLON DE TAILLE n
ET DE MOYENNE \bar{x} .



LUI ET NOUS SOMMES SÛRS À 95 % QUE LA MOYENNE μ EST DANS L'INTERVALLE
 $\bar{x} \pm 1,96s(\bar{x})$.



COMME AUPARAVANT, POUR UN SEUIL
DE CONFIANCE $1 - \alpha$, NOUS REMPLAÇONS
1,96 PAR $Z_{\alpha/2}$.



REVISITONS LES DONNÉES DES ÉTUDIANTS DU CHAPITRE 2. SUPPOSONS QUE NOS $n = 92$ ÉTUDIANTS REPRÉSENTENT UN ÉCHANTILLON ALÉATOIRE SIMPLE DE L'ENSEMBLE DES ÉTUDIANTS DE PENN STATE.



LA MOYENNE D'ÉCHANTILLON \bar{X} ÉTAIT DE 145,20 LIVRES ET L'ÉCART-TYPE ÉCHANTILLON s ÉTAIT DE 23,7. CELA FAIT UNE ERREUR-TYPE ESTIMÉE DE :

$$s(\bar{X}) = \frac{23,7}{\sqrt{92}} = 2,47$$

ET NOUS SOMMES CONFIANTS À 95 % QUE LE POIDS MOYEN DE **TOUS** LES ÉTUDIANTS DE PENN STATE EST DANS L'INTERVALLE :

$$\begin{aligned} \bar{x} \pm 1,96s(\bar{X}) \\ = 145,2 \pm (1,96)(2,47) \\ = 145,2 \pm 4,80 \text{ LIVRES } (65,9 \pm 2,2 \text{ kg}) \end{aligned}$$

EN RÉSUMÉ : POUR UN ÉCHANTILLON ALÉATOIRE SIMPLE DE GRANDE TAILLE, L'INTERVALLE DE CONFIANCE À $1 - \alpha$ EST :

MOYENNE DE POPULATION, μ

$$\mu \in \bar{x} \pm z_{\alpha/2}s(\bar{X})$$

où

$$s(\bar{X}) = \frac{s}{\sqrt{n}}$$

PROPORTION DE POPULATION, p

$$p \in \hat{p} \pm z_{\alpha/2}s(\hat{p})$$

où

$$s(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

LA LARGEUR DES DEUX INTERVALLES EST CONTRÔLÉE PAR LE SEUIL DE CONFIANCE $1 - \alpha$ ET LA TAILLE D'ÉCHANTILLON n .

MAINTENANT SÉNATEUR, QUE DIRIEZ-VOUS D'UN TRAVAIL DANS MON ENTREPRISE DE SONDAGE ?



t de Student (encore!)

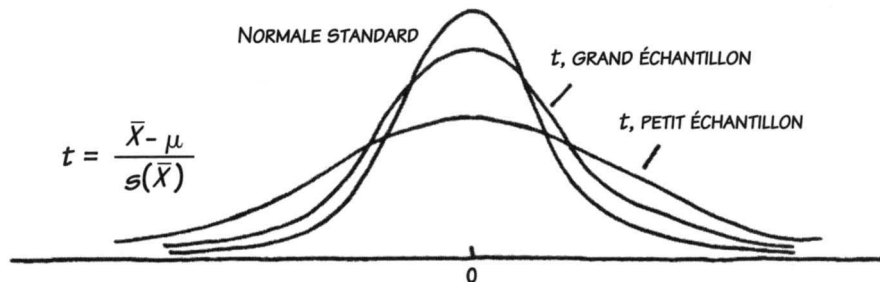
COMME NOUS L'AVONS VU DANS LE CHAPITRE 6, LA STATISTIQUE

$$\frac{\bar{X} - \mu}{s(\bar{X})}$$

EST PRESQUE NORMALE LORSQUE LA TAILLE D'ÉCHANTILLON EST GRANDE. POUR DES PETITS ÉCHANTILLONS ($n = 5, 10, 25, \dots$) CE N'EST PLUS LE CAS ET NOUS DEVONS UTILISER LE t DE STUDENT.



REGARDONS LE t DE STUDENT D'UN PEU PLUS PRÈS. NOUS AVONS MENTIONNÉ QUE SA DISTRIBUTION EST PLUS ÉTALÉE QU'UNE DISTRIBUTION NORMALE, ET QUE LA DISPERSION DÉPEND DE LA TAILLE D'ÉCHANTILLON.



LA DÉCOUVERTE DE GOSSET FUT DE **QUANTIFIER** CETTE RELATION. SI n EST LA TAILLE D'ÉCHANTILLON, IL DÉFINIT $n - 1$ COMME LE NOMBRE DE **degrés de liberté** DE L'ÉCHANTILLON.

L'IDÉE GÉNÉRALE : SOIT n DES OBSERVATIONS DONNÉES x_1, x_2, \dots, x_n . ON UTILISE UN DEGRÉ POUR LE CALCUL DE \bar{x} , LAISSANT $n - 1$ ÉLÉMENTS INDÉPENDANTS D'INFORMATION.

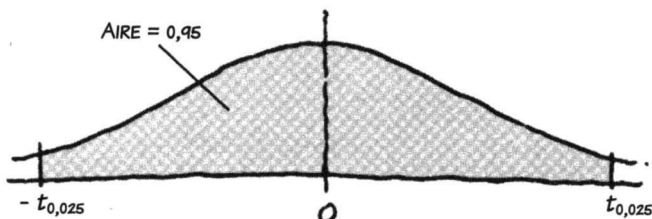


GOSSET CALCULA DES TABLES DE LA DISTRIBUTION t POUR DIFFÉRENTES TAILLES D'ÉCHANTILLON, ET DONC DIFFÉRENTS DEGRÉS DE LIBERTÉ. RAPPELONS QUE PLUS LES **DEGRÉS DE LIBERTÉ AUGMENTENT**, PLUS LE t DE STUDENT TEND VERS UNE LOI NORMALE STANDARD.



CONNAISSANT LA TAILLE D'ÉCHANTILLON n , ON PREND LE t DE STUDENT AVEC $n - 1$ DEGRÉS DE LIBERTÉ.

COMME POUR LA LOI NORMALE, NOUS AVONS UN SEUIL DE CONFIANCE À 95 % EN CHERCHANT LE FRACTILE $t_{0,025}$ AU-DELÀ DUQUEL L'AIRE SOUS LA COURBE EST DE 0,025.



COMME LA COURBE EST PLUS APLATIE QU'UNE NORMALE, $t_{0,025}$ EST PLUS LOIN QUE $z_{0,025}$.



POUR UN INTERVALLE DE CONFIANCE À $(1 - \alpha)$, ON TROUVE LE FRACTILE $t_{\alpha/2}$ TEL QUE $P(t \geq t_{\alpha/2}) = \alpha/2$. VOICI UNE TABLE RÉDUITE DE FRACTILES POUR LE t DE STUDENT.

$1 - \alpha$		0,8	0,9	0,95	0,99
α		0,2	0,1	0,05	0,01
$\alpha/2$		0,1	0,05	0,025	0,005
DEGRÉS DE LIBERTÉ	1	3,08	6,31	12,71	63,66
	10	1,37	1,81	2,23	3,17
	30	1,31	1,70	2,04	2,75
	100	1,29	1,66	1,98	2,63
	∞	1,28	1,64	1,96	2,58

CHAQUE COLONNE REPRÉSENTE UN SEUIL DE CONFIANCE FIXÉ EN FONCTION D'UN NOMBRE CROISSANT DE DEGRÉS DE LIBERTÉ. PLUS LE NOMBRE DE DEGRÉS DE LIBERTÉ AUGMENTE, PLUS LE FRACTILE TEND VERS $z_{\alpha/2}$, LE FRACTILE DE LA LOI NORMALE.

LA LARGEUR DE NOTRE INTERVALLE DE CONFIANCE SE DÉDUIT DE LA DÉFINITION DE t .

$$t = \frac{\bar{X} - \mu}{s(\bar{X})}$$

DONC AU SEUIL DE CONFIANCE $(1 - \alpha)$,

$$1 - \alpha = P(\bar{x} - t_{\alpha/2}s(\bar{X}) \leq \mu \leq \bar{x} + t_{\alpha/2}s(\bar{X}))$$

NOTE :
C'EST EXACTEMENT
COMME POUR LES
GRANDS ÉCHANTILLONS
MAIS AVEC t
AU LIEU DE z !



NOUS EN DÉDUISONS : SI L'ÉCHANTILLON EST DE TAILLE n ET DE MOYENNE \bar{x} , ALORS NOUS SOMMES SÛRS À $(1 - \alpha)$ QUE LA MOYENNE DE POPULATION μ APPARTIENT À L'INTERVALLE :

$$\mu \in \bar{x} \pm t_{\alpha/2}s(\bar{X})$$

OÙ $s(\bar{X}) = s/\sqrt{n}$, ET $t_{\alpha/2}$ EST LE FRACTILE DU t DE STUDENT AVEC $n - 1$ DEGRÉS DE LIBERTÉ.



MÉMORISE
BIEN CELA...

ENCORE
ÉVEILLÉ?



NOTE :

POUR ÊTRE PLUS RIGOUREUX, L'UTILISATION DU t DE STUDENT DÉPENDAIT DE L'HYPOTHÈSE QUE LA DISTRIBUTION DE LA POPULATION ÉTAIT NORMALE. DANS LA PRATIQUE, LES INTERVALLES DE CONFIANCE BASÉS SUR LE t DE STUDENT RESTENT RELATIVEMENT VALABLES, MÊME SI LA POPULATION N'EST QU'APPROXIMATIVEMENT EN FORME DE CLOCHE.

Exemple : L'ENTREPRISE **CAMÉLÉON AUTOMOBILES** DOIT EFFECTUER DES CRASH-TESTS SUR SES VOITURES AFIN D'ÉVALUER LE COÛT MOYEN DE RÉPARATION D'UNE COLLISION FRONTALE À 16 km/h. C'EST TRÈS COÛTEUX! ELLE DÉCIDE DONC DE FAIRE UN TEST SUR SEULEMENT CINQ CAMÉLÉON.



ON OBTIENT LES DONNÉES DE DOMMAGES : 150 €, 400 €, 720 €, 500 € ET 930 €.

LA MOYENNE D'ÉCHANTILLON EST :

$$\bar{x} = 540 \text{ €}$$

L'ÉCART-TYPE EST :

$$s = 299 \text{ €}$$

HUM! ÇA AMÉLIORE
LE DESIGN.

ON PEUT VÉRIFIER s À LA CALCULATRICE. C'EST :

$$\sqrt{\frac{1}{4}((150 - 540)^2 + (400 - 540)^2 + (720 - 540)^2 + (500 - 540)^2 + (930 - 540)^2)}$$

OÙ POUVONS-NOUS DONC PLACER LA MOYENNE AVEC 95 % DE CONFIANCE ?

NOUS TROUVONS NOTRE FRACTILE $t_{0,025}$ AVEC 4 DEGRÉS DE LIBERTÉ, QUI VAUT 2,78.

		$1 - \alpha$	0,8	0,9	0,95	0,99
		α	0,2	0,1	0,05	0,01
		$\alpha/2$	0,1	0,05	0,025	0,005
DEGRÉS DE LIBERTÉ	1		3,08	6,31	12,71	63,66
	2		1,89	2,92	4,3	9,92
	3		1,64	2,35	3,18	5,84
	4		1,53	2,13	2,78	4,6
	5		1,48	2,02	2,57	4,03

ET ON LE PLONGE DANS :

$$\mu \in \bar{x} \pm 2,78 \frac{s}{\sqrt{n}}$$

$$\mu \in 540 \pm 2,78 \frac{299}{\sqrt{5}}$$

$$\mu \in \mathbf{540 \pm 372}$$



TOUT CE QU'ON PEUT AFFIRMER AVEC 95 % DE CONFIANCE, C'EST QUE LES DOMMAGES MOYENS SERONT ENTRE 168 ET 912 €.



MAIS JE SUIS CONFIANT
À 0 % QUE CELA COÛTERA
EXACTEMENT 3,49 €...

L'ENTREPRISE PEUT SOIT
S'EN SATISFAIRE, SOIT FAIRE
D'AUTRES TESTS.

POUR CALCULER CET INTERVALLE DE CONFIANCE EN UTILISANT LE t DE STUDENT, NOUS AVONS FAIT UNE **HYPOTHÈSE IMPLICITE**. NOUS AVONS SUPPOSÉ QUE LES COÛTS DE RÉPARATION ÉTAIENT APPROXIMATIVEMENT **DISTRIBUÉS NORMALEMENT**, AINSI SI ON CRASHE 1000 CAMÉLÉON, L'HISTOGRAMME DES DOMMAGES DOIT ÊTRE SYMÉTRIQUE ET EN FORME DE MONTICULE. ON NE PEUT PAS DÉDUIRE CELA DE CET ÉCHANTILLON DE SEULEMENT 5 VOITURES... MAIS PEUT-ÊTRE QUE DES ANNÉES D'EXPÉRIENCE AVEC D'ANCIENS MODÈLES ONT DÉJÀ FOURNI DES HISTOGRAMMES NORMALEMENT DISTRIBUÉS DE DOMMAGE SUR LES PARTIES AVANT. CE TYPE D'INFORMATION JUSTIFIERAIT NOTRE UTILISATION DU t DE STUDENT.



ET POUR L'ARRIÈRE ?

LA QUEUE REPOUSSE
D'ELLE-MÊME.
C'EST UNE DES OPTIONS
DES CAMÉLÉON.

EN RÉSUMÉ (!),
 NOUS DISPOSONS
 MAINTENANT DE TROIS
 MÉTHODES POUR TROUVER
 DES INTERVALLES
 DE CONFIANCE. POUR
 DES PROPORTIONS COMME
 POUR DES MOYENNES
 AVEC DE GRANDS
 ÉCHANTILLONS,
 NOUS CHERCHONS $z_{\alpha/2}$
 DANS UNE TABLE NORMALE.
 POUR DES MOYENNES
 DE PETITS ÉCHANTILLONS,
 DISONS $n \leq 30$,
 ON TROUVE $t_{\alpha/2}$
 DANS LA TABLE
 DE STUDENT.



DANS TOUS LES CAS, LA LARGEUR DE L'INTERVALLE EST CE FRACTILE MULTIPLIÉ
 PAR L'ERREUR-TYPE ESTIMÉE.

$$z_{\alpha/2} s(\hat{p})$$

$$z_{\alpha/2} s(\bar{x})$$

$$t_{\alpha/2} s(\bar{x})$$

ET TOUTES CES ERREURS-TYPES SONT PROPORTIONNELLES AU NOMBRE MAGIQUE.

$$\frac{1}{\sqrt{n}}$$

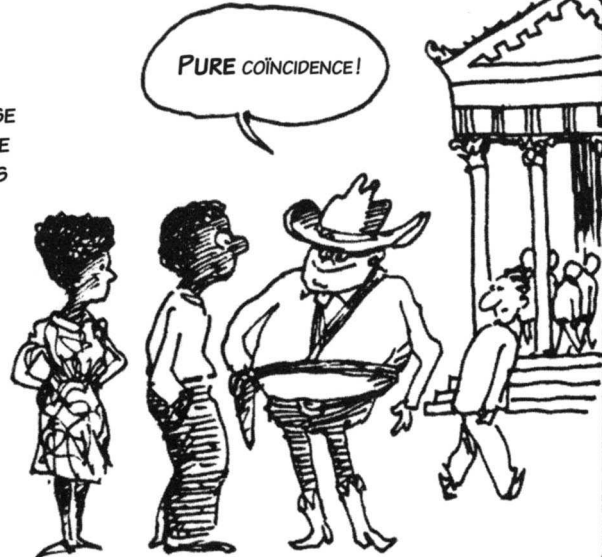
Chapitre 8

Tests d'hypothèses

NOUS ENTRONS MAINTENANT DANS UN NOUVEAU DOMAINE.
LES GOUVERNEMENTS, LES ENTREPRISES ET LES SCIENCES DURES
ET MOLLES UTILISENT TOUS ET ABUSENT SOUVENT DE CES TESTS
DE SIGNIFICATION. IL S'AGIT DE RÉPONDRE À LA QUESTION :
« **CES OBSERVATIONS ONT-ELLES RÉELLEMENT EU LIEU
PAR HASARD ?** »



NOUS COMMENÇONS
AVEC UN EXEMPLE PROVENANT
DU DOMAINE JURIDIQUE : PLUSIEURS
JUGEMENTS ONT ÉTÉ REMIS EN CAUSE
DANS LE SUD DES ÉTATS-UNIS ENTRE
1960 ET 1980. DES TÉMOINS EXPERTS
PRÉSENTÈRENT DES CAS DE **BIAIS**
RACIAUX DANS LA **SÉLECTION**
DES JURYS.



LES PANELS DE JURÉS SONT THÉORIQUEMENT COMPOSÉS ALÉATOIREMENT À PARTIR D'UNE LISTE
DE CITOYENS ÉLIGIBLES. CEPENDANT, DANS LES ÉTATS DU SUD DES ANNÉES 1950 ET 1960, IL Y AVAIT
PEU D'AFRO-AMÉRICAINS DANS CES JURYS. CERTAINS AVOCATS CONTESTÈRENT DONC LE VERDICT.
EN APPEL, UN TÉMOIN EXPERT EN STATISTIQUES FOURNIT CETTE PREUVE :

1) 50 % DE CITOYENS ÉLIGIBLES
ÉTAIENT AFRO-AMÉRICAINS.



2) SUR UN PANEL DE 80 JURÉS
POTENTIELS, SEULEMENT **QUATRE**
ÉTAIENT AFRO-AMÉRICAINS.



CELA POUVAIT-IL ÊTRE LE RÉSULTAT
D'UN **PUR HASARD** ?

POUR SIMPLIFIER LA DISCUSSION, SUPPOSONS QUE LA SÉLECTION PARMI LES JURÉS POTENTIELS SOIT **ALÉATOIRE**. ALORS, LE NOMBRE D'AFRO-AMÉRICAINS SUR UN PANEL DE 80 PERSONNES SERA UNE **VARIABLE ALÉATOIRE X** AVEC $n = 80$ TIRAGES ET $p = 1/2$.



AINSI, LA PROBABILITÉ D'OBTENIR UN JURY DE MOINS DE QUATRE AFRO-AMÉRICAINS EST $P(X \leq 4)$, QUI S'AVÈRE ÊTRE DE L'ORDRE DE 0,00000000000000000014 (!).



COMME LA PROBABILITÉ EST INFINITÉSIMALE, LE PANEL EN QUESTION AVEC SEULEMENT QUATRE NOIRS AMÉRICAINS FOURNIT UNE **PREUVE SOLIDE CONTRE L'HYPOTHÈSE** D'UNE SÉLECTION ALÉATOIRE.



POUR PORTER L'ESTOCADÉ, LE STATISTICIEN FAIT REMARQUER QUE CETTE PROBABILITÉ EST PLUS FAIBLE QUE CELLE D'OBTENIR **TROIS QUINTES FLUSH ROYALES CONSÉCUTIVES** AU POKER.



LE JUGE REJETTE DONC L'HYPOTHÈSE D'UNE SÉLECTION ALÉATOIRE.



SUIVONS LE MÊME PROCÉDÉ QUE PRÉCÉDEMMENT
POUR DÉVELOPPER LES **QUATRE ÉTAPES**
FORMELLES D'UN TEST D'HYPOTHÈSE STATISTIQUE.

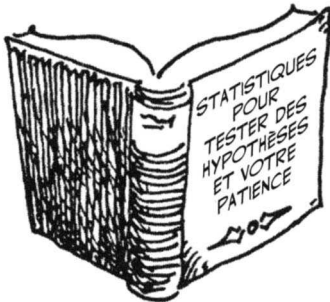
**Étape 1 : FORMULER TOUTES
LES HYPOTHÈSES.**

H_0 L'HYPOTHÈSE NULLE EST
GÉNÉRALEMENT QUE LES OBSERVATIONS
SONT STRICTEMENT LE RÉSULTAT
DU **HASARD**.

H_a L'HYPOTHÈSE ALTERNATIVE
INDIQUE QU'IL Y A UN EFFET RÉEL,
QUE LES OBSERVATIONS SONT LE RÉSULTAT
DE CET EFFET ET D'UNE COMPOSANTE
ALÉATOIRE.



Étape 2 : LE TEST STATISTIQUE.
IDENTIFIER LA STATISTIQUE QUI FOURNIRA
DES PREUVES CONTRE L'HYPOTHÈSE NULLE.



DANS L'AFFAIRE DE JUSTICE, H_0
STIPULE QUE LE JURY ÉTAIT **CHOISI**
AU HASARD DANS LA POPULATION
GLOBALE. LES AFRO-AMÉRICAINS
ONT UNE PROBABILITÉ $p = 0,50$
D'ÊTRE CHOISIS.

H_a STIPULE QUE LES AFRO-AMÉRICAINS
SONT **MOINS SÉLECTIONNÉS**
DANS LES JURYS QUE LEUR PROPORTION
DANS LA POPULATION : $p < 0,50$.



ICI, LA STATISTIQUE DE TEST EST
LA VARIABLE ALÉATOIRE **BINOMIALE** X
AVEC $p = 0,50$ ET $n = 80$.



Étape 3 : VALEUR P : UNE CONSTATATION PROBABILISTE POUR RÉPONDRE À LA QUESTION : SI L'HYPOTHÈSE NULLE EST VRAIE, ALORS QUELLE EST LA PROBABILITÉ D'OBSERVER UNE STATISTIQUE DE TEST AU MOINS AUSSI EXTRÊME QUE CELLE OBSERVÉE ?

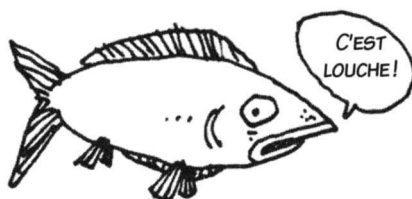


Étape 4 : COMPARER LA VALEUR P À UN SEUIL DE SIGNIFICATION, α .

α AGIT COMME UN SEUIL LIMITE EN DESSOUS DUQUEL NOUS NOUS ACCORDONS À DIRE QUE L'EFFET EST STATISTIQUEMENT SIGNIFICATIF. AUTREMENT DIT, SI

$$P \leq \alpha$$

ALORS NOUS REJETONS L'HYPOTHÈSE NULLE H_0 ET DÉCLARONS QUE QUELQUE CHOSE D'AUTRE A EU LIEU.



DANS NOTRE EXEMPLE, LA VALEUR P VAUT

$$P(X \leq 4 \mid p = 0,50 \text{ ET } n = 80)$$

$$= 1,4 \times 10^{-18}$$

ON A CALCULÉ LA VALEUR P AVEC LES MOYENS MODERNES EN UTILISANT UN LOGICIEL DE STATISTIQUES.



α EST SOUVENT PRIS À 0,05 OU 0,01. DANS LE CAS DU JURY, ET POUR L'EFFET DRAMATIQUE, LE STATISTICIEN A FAIT ALLUSION À UN α ABSURDEMENT FAIBLE DE $3,6 \times 10^{-18}$, QUI REPRÉSENTE LES CHANCES D'AVOIR TROIS QUINTES FLUSH ROYALES DE SUITE.



DANS LES TRAVAUX SCIENTIFIQUES, ON UTILISE SOUVENT UN SEUIL FIXE DE α ÉGAL À 0,05 OU 0,01. CES VALEURS SONT UN VESTIGE DE L'ÈRE PRÉ-ORDINATEUR OÙ L'ON DEVAIT SE RÉFÉRER À DES TABLES IMPRIMÉES POUR TROUVER DES VALEURS CRITIQUES PRÉSÉLECTIONNÉES. POURTANT DE NOS JOURS, LES JOURNAUX SCIENTIFIQUES PUBLIENT SEULEMENT DES RÉSULTATS LORSQUE LA VALEUR P EST $\leq 0,05$.



DANS LES COURS DE JUSTICE,
LE STANDARD EST PLUS FLEXIBLE...

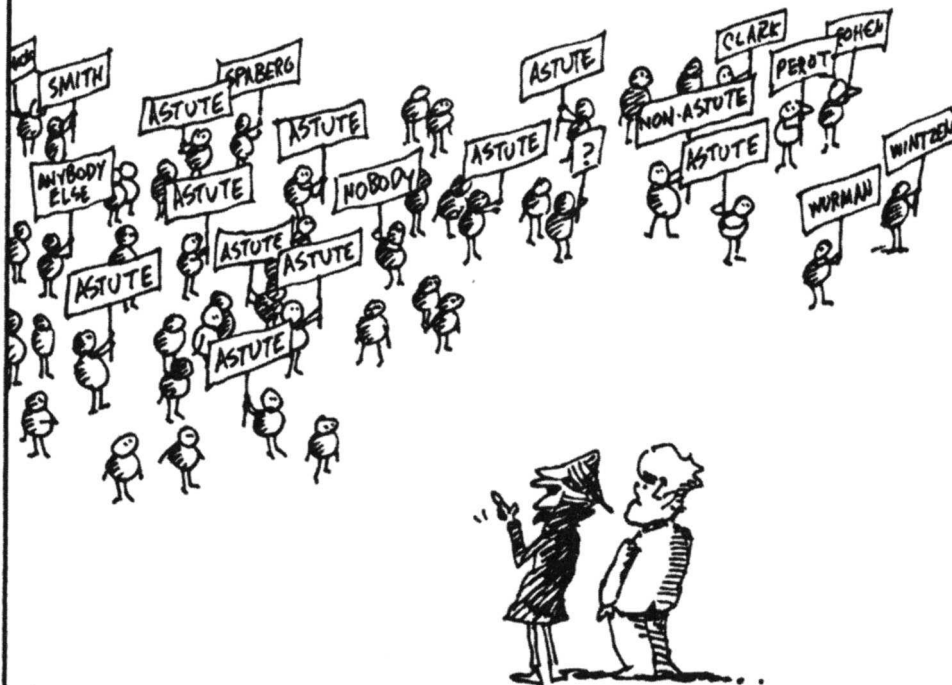


Test de signification de proportion pour de GRANDS ÉCHANTILLONS

L'EXEMPLE DU JURY EST UN CAS PARTICULIER
D'UN PROBLÈME PLUS GÉNÉRAL OÙ L'HYPOTHÈSE NULLE
EST DE LA FORME $p = p_0$ ET OÙ p_0 EST UNE CERTAINE
PROBABILITÉ (DANS NOTRE CAS 0,5). MAINTENANT,
ÉTUDIONS LE CAS GÉNÉRAL.
TESTONS L'HYPOTHÈSE $p = p_0$.



COMME D'HABITUDE, NOUS SUPPOSONS QUE NOUS AVONS UNE ÉNORME POPULATION...
NOUS OBSERVONS UN GRAND ÉCHANTILLON... ET NOUS TROUVONS QU'UNE PROPORTION
DE L'ÉCHANTILLON VÉRIFIE UNE PROPRIÉTÉ.



À PARTIR DE CETTE OBSERVATION, NOUS VOULONS SAVOIR SI LA VRAIE PROPORTION
DE POPULATION EST (PAR EXEMPLE) PLUS GRANDE QU'UNE VALEUR p_0 . PAR EXEMPLE,
LE SÉNATEUR ASTUTE QUI A TROUVÉ $\hat{p} = 0,55$ VEUT SAVOIR SI $p > 0,5$, CE QUI LUI
GARANTIRAIT LA MAJORITÉ.

Étape 1

L'HYPOTHÈSE NULLE EST :

$$H_0 : p = p_0$$

L'HYPOTHÈSE ALTERNATIVE DÉPEND DE LA DIRECTION DE L'EFFET CONSIDÉRÉ. DANS LE CAS DU SÉNATEUR ASTUTE :

$$H_a : p > p_0$$

MAIS DANS D'AUTRES CAS, L'HYPOTHÈSE ALTERNATIVE PEUT ÊTRE :

$$H_a : p < p_0$$

OU

$$H_a : p \neq p_0$$

PAR EXEMPLE, DANS LE CAS DE LA SÉLECTION DU JURY, L'HYPOTHÈSE ALTERNATIVE ÉTAIT :

$$H_a : p < 0,5$$

MAIS D'AUTRES FOIS, NOUS VOULONS SAVOIR SI p EST DIFFÉRENT OU NON D'UNE VALEUR. PAR EXEMPLE, SI ON TESTE L'HONNÊTÉTÉ D'UN LANCER DE PIÈCE, L'HYPOTHÈSE ALTERNATIVE EST :

$$H_a : p \neq 0,5$$

MAIS NOUS N'AVONS AUCUNE IDÉE A PRIORI SI LE BIAIS EST SUR LE CÔTÉ FACE OU PILE.



Étape 2

LA STATISTIQUE DE TEST EST :

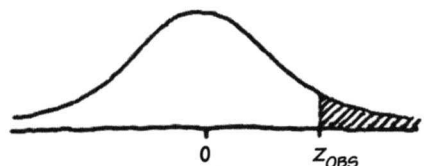
$$z_{OBS} = \frac{(\hat{p} - p_0)}{\sqrt{p_0(1 - p_0)/n}}$$

QUI MESURE L'ÉLOIGNEMENT DE \hat{p} PAR RAPPORT À p_0 . AVEC L'HYPOTHÈSE NULLE, z_{OBS} A UNE DISTRIBUTION NORMALE STANDARD.

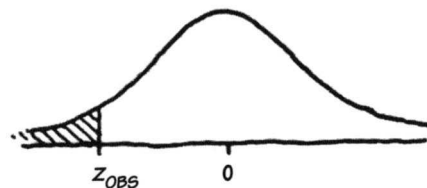
Étape 3

LA VALEUR P DÉPEND DU TYPE DE L'HYPOTHÈSE ALTERNATIVE :

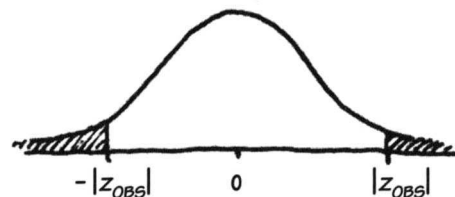
A) TEST UNILATÉRAL À DROITE, $H_a : p > p_0$.
UTILISE LA VALEUR P DÉFINIE PAR $P(z > z_{OBS})$.



B) TEST UNILATÉRAL À GAUCHE, $H_a : p < p_0$.
UTILISE LA VALEUR P DÉFINIE PAR $P(z < z_{OBS})$.



C) TEST BILATÉRAL À GAUCHE, $H_a : p \neq p_0$.
UTILISE LA VALEUR P DÉFINIE PAR $P(|z| > |z_{OBS}|)$.



DANS LE CAS DU SÉNATEUR ASTUTE :

1) LES HYPOTHÈSES SONT :

$$H_0 : p = 0,5$$

$$H_a : p > 0,5$$

2) LA STATISTIQUE DE TEST EST :

$$z_{OBS} = \frac{(0,55 - 0,50)}{\sqrt{(0,50)(0,50)/1000}} = 3,16$$

3) LA VALEUR P EST :

$$P(z \geq z_{OBS}) = P(z \geq 3,16) = 0,0008$$

4) ÉTANT CONSERVATEUR, ASTUTE PREND UN SEUIL DE SIGNIFICATION α DE 0,01 ET IL OBSERVE QUE :

$$P(z > z_{OBS}) = 0,0008 < \alpha$$

LE SÉNATEUR REJETTE DONC L'HYPOTHÈSE NULLE, ET LUI ET SON ÉQUIPE PEUVENT ÊTRE SÛRS D'ÊTRE EN TÊTE.

VOUS POUVEZ CONTRIBUER
MAINTENANT...



NDT : LORSQUE L'HYPOTHÈSE ALTERNATIVE EST UNE INÉGALITÉ (DU TYPE $H_0 : p < p_0$), CERTAINS STATISTICIENS PRENNENT COMME HYPOTHÈSE NULLE LE CONTRAIRE DE L'HYPOTHÈSE ALTERNATIVE (SOIT $H_0 : p \geq p_0$). CELA NE CHANGE EN RIEN L'ANALYSE CAR C'EST L'HYPOTHÈSE ALTERNATIVE QUI DÉFINIT LE CRITÈRE DE REJET OU NON.

Test de **MOYENNE** de population pour de **GRANDS ÉCHANTILLONS**

VOICI MAINTENANT UN TEST DE SIGNIFICATION
QUI PEUT ÊTRE UTILISÉ POUR UN **CONTRÔLE**
PAR ÉCHANTILLONNAGE, UNE APPLICATION
INDUSTRIELLE IMPORTANTE.



MAIS BIEN SÛR L'ÉPICERIE N'A AUCUNE INTENTION DE PESER CHACUNE DES BOÎTES EXPÉDIÉES. SES EMPLOYÉS SE CONTENTERONT D'UTILISER LES STATISTIQUES !



TOUT D'ABORD, ILS CHOISSENT LEURS HYPOTHÈSES :

$$H_0 : \mu = 16 \text{ OZ.}$$

$$H_a : \mu < 16 \text{ OZ.}$$

REJETER L'HYPOTHÈSE IMPLIQUE DE REFUSER LA LIVRAISON DES CÉRÉALES.



ENSUITE, ILS CHOISSENT UNE STATISTIQUE DE TEST. MAINTENANT, CELA DEVRAIT ÊTRE COMME UN RÉFLEXE D'EXTENSION DU GENOU DE DIRE QUE L'ÉCART DE L'ÉCHANTILLON À LA MOYENNE EST :

$$\frac{\bar{X} - \mu_0}{s(\bar{X})} = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$$

OÙ s EST L'ÉCART-TYPE DE L'ÉCHANTILLON. SOUS L'HYPOTHÈSE NULLE, CELA SUIT PRESQUE UNE LOI NORMALE STANDARD AVEC DE GRANDS ÉCHANTILLONS, D'APRÈS LE THÉORÈME CENTRAL LIMITE.



EN SAUTANT L'ÉTAPE 3 POUR LE MOMENT, ILS FIXENT UN SEUIL DE SIGNIFICATION. AYANT TOUS RATÉ LA PLACE DE MAJOR EN SCIENCES, NOS ÉPICIERS PENSENT QUE $\alpha = 0,05$ SONNE BIEN.



ILS PRENNENT UN ÉCHANTILLON
ALÉATOIRE SIMPLE DE 49 BOÎTES,
LES PÈSENT ET DÉTERMINENT
LES STATISTIQUES DE L'ÉCHANTILLON.

$$\bar{x} = 15,9 \text{ OZ}$$

$$s = 0,35 \text{ OZ}$$

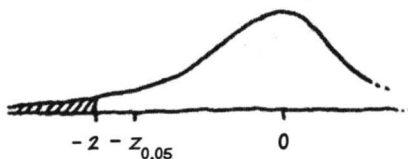
C'EST UN PEU LÉGER,
MAIS EST-CE SIGNIFICATIF?



ILS RENTRENT LES VALEURS DANS LA STATISTIQUE DE TEST
POUR TROUVER :

$$Z_{OBS} = \frac{15,9 - 16}{0,35/\sqrt{49}} = -2$$

MAINTENANT, ILS CALCULENT LEUR
VALEUR P : $P(z < -2 \mid H_0) = 0,0227$



COMME C'EST MOINS QUE LE SEUIL
DE SIGNIFICATION DE 0,05, L'ÉPICERIE
AUTHENTIQUE REJETTE L'HYPOTHÈSE
NULLE, ET LA MARCHANDISE.

TU PEUX
LA RAMENER,
ESPÈCE D'ARTISTE
RATÉ!



QUE S'EST-IL
PASSÉ?



J'AVAIS LES CROCS, MAN...
JE NE PENSais PAS QUE QUELQU'UN
REMARQUERAIT SI J'EN MANGEAIS
UN PEU DANS CHAQUE BOÎTE...

Test de MOYENNE de population pour de PETITS ÉCHANTILLONS



NOUS REVENONS DANS L'ENTREPRISE CAMÉLÉON AUTOMOBILES ET SES CRASH-TESTS. LA **COMPAGNIE D'ASSURANCES INTÈGRE** ASSURERA UNE AUTO SEULEMENT SI LA MOYENNE DES COÛTS DE RÉPARATION APRÈS UNE COLLISION À 16 km/h S'ÉLÈVE À MOINS DE 1000 €. LA COMPAGNIE UTILISE UN SEUIL DE SIGNIFICATION STANDARD $\alpha = 0,05$. AINSI :

$H_0 : \mu \geq 1000 \text{ €}$ LE COÛT MOYEN EST TROP ÉLEVÉ.

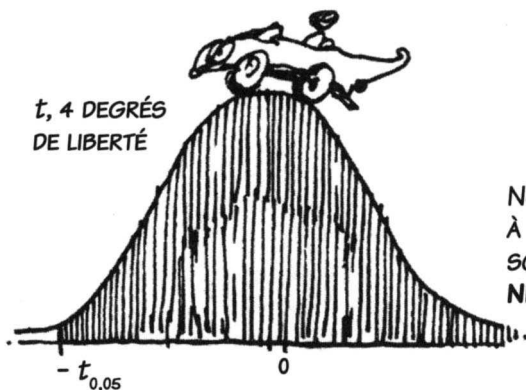
$H_a : \mu < 1000 \text{ €}$ LE COÛT MOYEN EST CORRECT.

LA STATISTIQUE DE TEST EST :

$$t = \frac{\bar{X} - \mu_0}{s(\bar{X})} \quad \text{OÙ } \mu_0 \text{ EST LA MOYENNE HYPOTHÉTIQUE DE 1000 €}.$$



SOUS L'HYPOTHÈSE NULLE, LA STATISTIQUE A UNE DISTRIBUTION t DE STUDENT AVEC 4 DEGRÉS DE LIBERTÉ.



NOUS VOULONS QUE LE t OBSERVÉ SOIT À GAUCHE DE $-t_{0,05}$ (CAR DE FAIBLES \bar{x} SONT SOUHAITÉS, $\bar{x} - \mu_0$ DOIT ÊTRE NÉGATIF POUR CONFIRMER H_a).

DEGRÉS
DE LIBERTÉ

	α		
	0,05	0,025	0,005
1	6,31	12,71	63,66
2	2,92	4,30	9,92
3	2,35	3,18	5,84
4	2,13	2,78	4,6
5	2,01	2,57	4,03

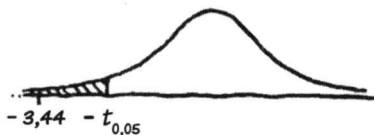
DANS LA TABLE DES VALEURS CRITIQUES t , ON VOIT QUE $t_{0,05} = 2,13$, ON REJETTE DONC H_0 SI :

$$t_{OBS} \leq -t_{0,05} = -2,13$$

NOUS SAVONS (VOIR CHAPITRE 7) QUE $\bar{x} = 540$ € ET $s = 299$ € POUR NOTRE PETIT ÉCHANTILLON DE 5 VOITURES, AINSI ON OBTIENT :

$$t_{OBS} = \frac{540 - 1000}{299/\sqrt{5}} = -3,44 < -t_{0,05}$$

FÉLICITATIONS!
PARLONS
MAINTENANT
DE VOS AUTRES
BESOINS EN
ASSURANCE...



LA VOITURE PASSE LE TEST... H_0 EST REJETÉE... ET LA POLICE D'ASSURANCE EST ÉMISE.

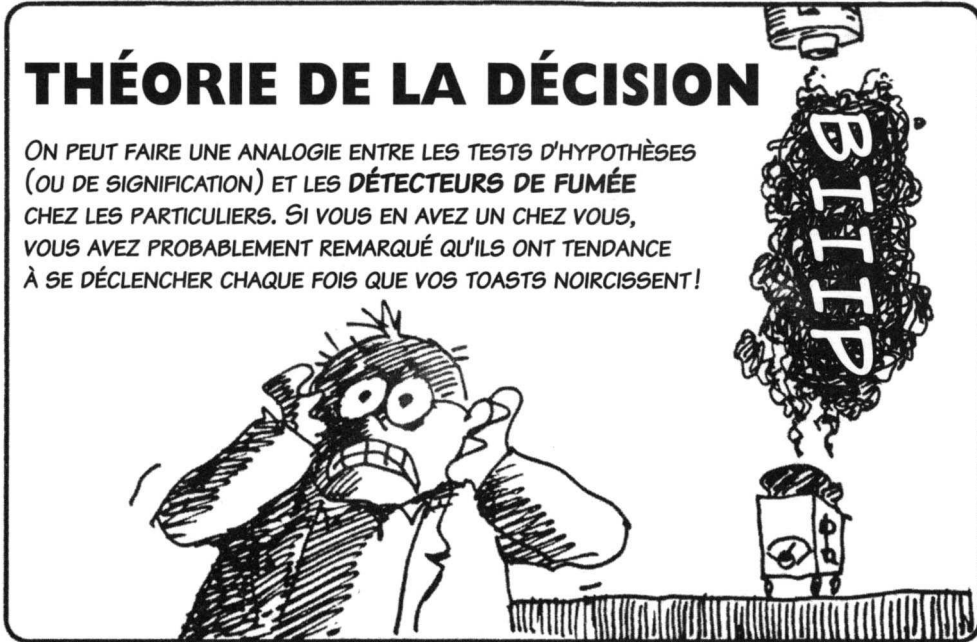


EST-CE QUE
CES MOUCHES ONT
UNE ASSURANCE-VIE?

C'EST UN EXEMPLE D'ÉCHANTILLONNAGE D'ACCEPTATION. L'HYPOTHÈSE NULLE EST QUE LES COÛTS DE RÉPARATION NE SONT PAS ACCEPTABLES. L'ENTREPRISE EST PRÉSUMÉE COUPABLE ET DOIT FOURNIR DES PREUVES SUFFISANTES DE SON INNOCENCE EN MONTRANT QUE LE PRODUIT EST CONFORME AUX SPÉCIFICATIONS.

THÉORIE DE LA DÉCISION

ON PEUT FAIRE UNE ANALOGIE ENTRE LES TESTS D'HYPOTHÈSES (OU DE SIGNIFICATION) ET LES **DÉTECTEURS DE FUMÉE** CHEZ LES PARTICULIERS. SI VOUS EN AVEZ UN CHEZ VOUS, VOUS AVEZ PROBABLEMENT REMARQUÉ QU'ILS ONT TENDANCE À SE DÉCLANCHER CHAQUE FOIS QUE VOS TOASTS NOIRCISSENT !



C'EST CE QUI S'APPELLE UNE **ERREUR DE TYPE I** : UNE ALARME SANS FEU. INVERSEMENT, UNE **ERREUR DE TYPE II** EST UN FEU SANS ALARME. TOUT CUISINIER SAIT ÉVITER LES ERREURS DE TYPE I : IL SUFFIT D'**ENLEVER LES PILES**. MALHEUREUSEMENT, CELA AUGMENTE LA FRÉQUENCE DES ERREURS DE TYPE II !



DE FAÇON SIMILAIRE, RÉDUIRE LES RISQUES D'ERREURS DE TYPE II, EN RENDANT, PAR EXEMPLE, L'ALARME HYPERSENSIBLE, PEUT ACCROÎTRE LE NOMBRE DE FAUSSES ALARMES.

ON PEUT RÉSUMER CES RÉSULTATS DANS UNE **TABLE DE DÉCISION** :

	PAS DE FEU	FEU
PAS D'ALARME	PAS D'ERREUR	TYPE II
ALARME	TYPE I	PAS D'ERREUR

MAINTENANT, VOYONS L'HYPOTHÈSE NULLE COMME ÉTANT LA CONDITION QU'IL N'Y A **PAS DE FEU**, ALORS QUE L'HYPOTHÈSE ALTERNATIVE SIGNIFIE QU'IL Y A LE FEU. L'ALARME CORRESPOND AU REJET DE L'HYPOTHÈSE NULLE.

	VÉRITABLE ÉTAT	
	H_0	H_a
ACCEPTER H_0	PAS D'ERREUR	TYPE II
REJETER H_0	TYPE I	PAS D'ERREUR

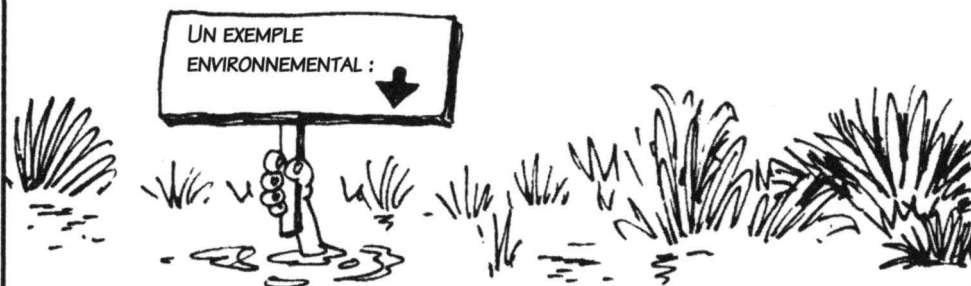
TOUS LES TESTS DE SIGNIFICATION QUE NOUS VENONS D'EFFECTUER DANS CE CHAPITRE METTENT L'ACCENT SUR LA PROBABILITÉ DE COMMETTRE UNE ERREUR DE TYPE I (C'EST-À-DIRE LA PROBABILITÉ QUE NOS OBSERVATIONS SE RÉALISENT LORSQUE H_0 EST VRAIE). NOUS AVONS POSÉ :

$$P(\text{REJETER } H_0 \mid H_0) = P(\text{ERREUR DE TYPE I} \mid H_0) = \alpha$$

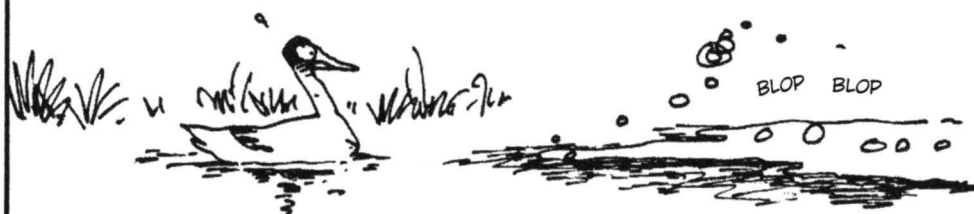
$1 - \alpha$ MESURE NOTRE CONFIANCE QUE QUAND L'ALARME SONNE, IL Y A EFFECTIVEMENT LE FEU. UNE CONFIANCE ÉLEVÉE SIGNIFIE QUE LE DÉCLENCHEMENT DE FAUSSES ALARMES EST RARE.



MAIS PARFOIS, CE QUE NOUS VOULONS VRAIMENT CONNAÎTRE C'EST LA PROBABILITÉ DE FAIRE UNE **ERREUR DE TYPE II**. EN D'AUTRES TERMES, QUELLE EST LA SENSIBILITÉ DE NOTRE «SYSTÈME D'ALARME» QUAND L'HYPOTHÈSE ALTERNATIVE EST VÉRIFIÉE ?



DANS LE PASSÉ, LES USINES QUI DÉVERSAIENT DES PRODUITS CHIMIQUES DANS LES COURS D'EAU DEVAIENT PROUVER QUE CETTE ACTION N'AVAIT AUCUN IMPACT SUR LA VIE AQUATIQUE EN AVAL. C'EST H_0 . LE POLLUEUR POUVAIT CONTINUER TANT QUE L'HYPOTHÈSE NULLE N'ÉTAIT PAS REJETÉE AU SEUIL DE SIGNIFICATION DE 5 %.



AINSI LE POLLUEUR QUI PENSAIT ÊTRE EN VIOLATION DES NORMES CONCEVAIT LUI-MÊME UN PROGRAMME DE CONTRÔLE DE POLLUTION INEFFICACE.



LE POLLUEUR EST RAVI PUISQUE, COMME NOTRE ALARME INCENDIE SANS PILES, SON TEST A TRÈS PEU OU AUCUNE CHANCE DE DÉCLENCHER L'ALARME.



FORMALISONS CETTE IDÉE. POUR DÉCRIRE LA **PROBABILITÉ D'UNE ERREUR DE TYPE II**, NOUS SORTONS UNE AUTRE LETTRE GRECQUE, BÊTA OU β .

$$\beta = P(\text{ACCEPTER } H_0 \mid H_a) \\ = P(\text{ERREUR DE TYPE II} \mid H_a)$$

LA **PUISSANCE** D'UN TEST EST DÉFINIE PAR $1 - \beta$. IL S'AGIT DE :

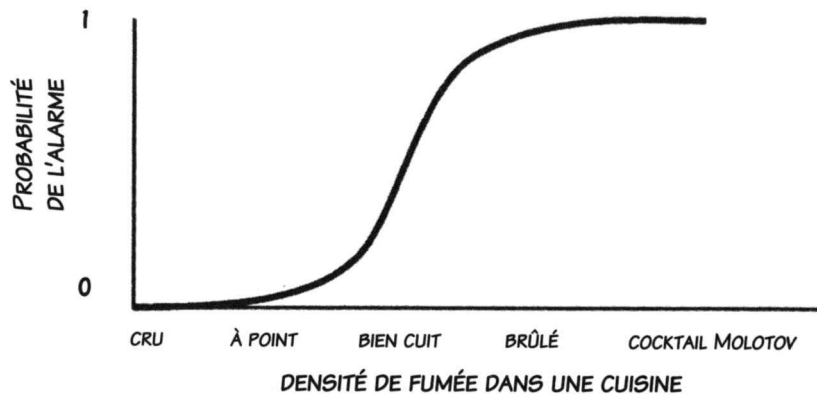
$$P(\text{REJETER } H_0 \mid H_a).$$



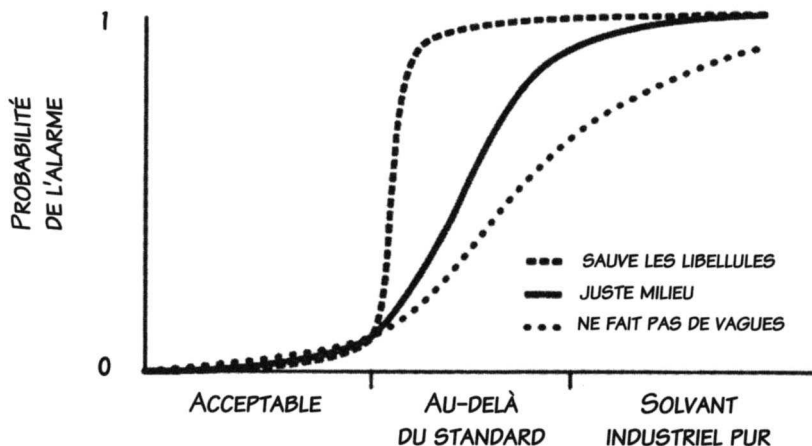
VOUS SEREZ CONTENT D'APPRENDRE QUE LES AGENCES ENVIRONNEMENTALES ONT ÉVOLUÉ VERS DES PROGRAMMES DE CONTRÔLE DE POLLUTION QUI DOIVENT MONTRER UNE PROBABILITÉ ÉLEVÉE DE DÉTECTER LES ÉVÉNEMENTS POLLUANTS IMPORTANTS. L'**ANALYSE DE PUISSANCE** REQUISE RÉVÈLE SOUVENT LES FAILLES CACHÉES DES PROGRAMMES DE CONTRÔLES.



UNE FAÇON DE VISUALISER L'EFFET DE LA PUISSANCE D'UN TEST EST DE TRACER LA PROBABILITÉ DE REJETER H_0 COMME UNE FONCTION DE L'ÉTAT DU SYSTÈME. DANS LE CAS DE L'ALARME INCENDIE, LA PROBABILITÉ CROÎT VERS 1 QUAND LA FUMÉE S'ÉPAISSIT.



POUR L'EXEMPLE SUR LA QUALITÉ DE L'EAU, L'AXE HORIZONTAL EST LA VRAIE CONCENTRATION DE POLLUANTS DANS L'EAU.



VOICI LES COURBES DE PUISSANCE DE TROIS PROGRAMMES DE CONTRÔLE. LE «**SAUVE LES LIBELLULES**» (COÛTE 5 MILLIONS DE €), LE «**JUSTE MILIEU**» (COÛTE 500 000 €), LE «**NE FAIT PAS DE VAGUES**» (COÛTE AUSSI 500 000 €, MAIS IL GARANTIT UN BON SPECTACLE!). PLUS LA PUISSANCE DU TEST EST IMPORTANTE, PLUS LA COURBE EST ABRUPTÉ.

FÉLICITATIONS !
AVEC CES SECTIONS QUI COUVRENT
LES BASES DES INTERVALLES DE CONFIANCE
ET DES TESTS D'HYPOTHÈSES, VOUS AVEZ TERMINÉ
VOTRE PREMIER COURS
DE STATISTIQUES GÉNÉRALES !

AH BON ? !



POURQUOI AVEZ-VOUS ALORS CETTE **SENSATION DE VIDE** DANS VOTRE ESTOMAC ?
C'EST PARCE QUE POUR UTILISER CES IDÉES DE FAÇON PRATIQUE, NOUS DEVONS
POUVOIR LES APPLIQUER À UNE VARIÉTÉ DE SITUATIONS QUE NOUS N'AVONS PAS
ENCORE VUES. C'EST PAR LÀ QUE NOUS ALLONS POURSUIVRE AVEC LA **COMPARAISON
DE DEUX POPULATIONS**.

OK!
AMENEZ LES POPULATIONS !



Chapitre 9

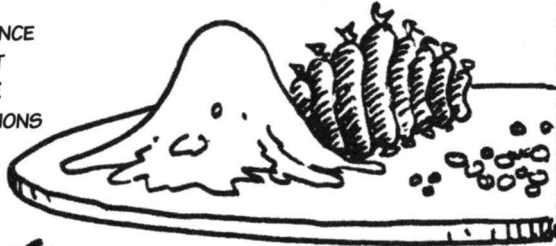
Comparaison de deux populations

*DANS LEQUEL NOUS APPRENNONS DE NOUVELLES RECETTES
EN UTILISANT DE VIEUX INGRÉDIENTS.*



LES DEUX DERNIERS CHAPITRES NOUS ONT EXPLIQUÉ LES INTERVALLES DE CONFIANCE ET LES TESTS D'HYPOTHÈSES EN UTILISANT LA **VIANDE ET LES POMMES DE TERRE** DES MODÈLES ALÉATOIRES : LES DISTRIBUTIONS NORMALES ET BINOMIALES.

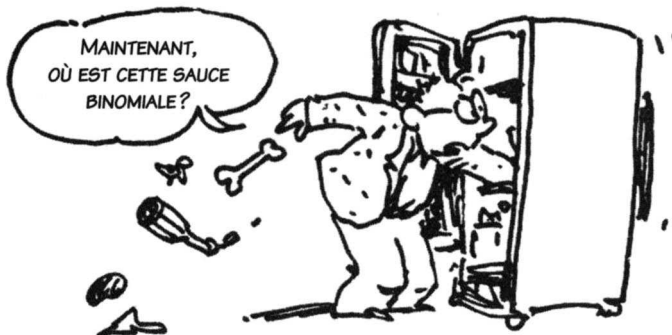
AVEC LA LOI NORMALE
DANS LE RÔLE DES PATATES!



MAIS CE QUI REND LES STATISTIQUES AUSSI STIMULANTES QUE LA CUISINE, C'EST LA DIVERSITÉ. COMME UN EXPERT EN CUISINE, LE STATISTICIEN PEUT « GOÛTER » LES INGRÉDIENTS D'UN PROBLÈME ET TROUVER LA FAÇON LA PLUS EFFICACE DE LES COMBINER DANS UNE RECETTE STATISTIQUE.



(LA RAISON POUR LAQUELLE LES LIVRES DE CUISINE ET LES MANUELS DE STATISTIQUES SONT SI LOURDS EST QU'ILS FOURNISSENT TOUS LES DEUX DES SOLUTIONS À UNE GRANDE VARIÉTÉ DE SITUATIONS!)



DANS CE CHAPITRE, NOUS ALLONS UTILISER NOS MÉTHODES « VIANDE-ET-PATATE » DANS DE NOUVELLES RECETTES QUI VONT NOUS AIDER À RÉPONDRE AUX QUESTIONS SUIVANTES :



EST-CE QUE PRENDRE DE L'ASPIRINE RÉGULIÈREMENT RÉDUIT LE RISQUE D'INFARCTUS ?



EST-CE QU'UN PESTICIDE PARTICULIER AUGMENTE LE RENDEMENT AGRICOLE ?



EST-CE QUE LES HOMMES ET LES FEMMES FAISANT LE MÊME TRAVAIL ONT DES SALAIRES DIFFÉRENTS ?



L'INGRÉDIENT COMMUN À CES QUESTIONS EST QU'IL EST POSSIBLE D'Y RÉPONDRE EN COMPARANT DEUX ÉCHANTILLONS ALÉATOIRES INDÉPENDANTS, UN POUR CHAQUE POPULATION.



PESTICIDE



SANS PESTICIDE

ET À LA FIN DU CHAPITRE, NOUS VERRONS UNE AUTRE MÉTHODE QUI PERMET DE COMPARER DEUX MOYENNES ET QUI NE NÉCESSITE PAS DE PRENDRE DEUX ÉCHANTILLONS ALÉATOIRES SIMPLES...



Comparaison de TAUX DE SUCCÈS (ou d'échecs) de deux populations

NOUS COMMENÇONS PAR UNE EXPÉRIENCE, EN PARTIE MENÉE PAR UNE ÉTUDE DE HARVARD, QUI CHERCHAIT À DÉTERMINER L'EFFICACITÉ DE L'**ASPIRINE** POUR RÉDUIRE LES CRISES CARDIAQUES. COMME DANS LA PLUPART DES ESSAIS CLINIQUES, LES RISQUES QU'UN INDIVIDU AIT LA MALADIE – ICI UNE CRISE CARDIAQUE – SONT TRÈS MINCES SUR UNE ANNÉE. MAIS NOUS VOULONS UNE RÉPONSE RAPIDE ! COMMENT FAIRE ?



LA SOLUTION SIMPLE MAIS COÛTEUSE EST DE TESTER UN GRAND NOMBRE D'INDIVIDUS EN PEU DE TEMPS. DANS CETTE ÉTUDE, 22 071 SUJETS (TOUS DES DOCTEURS VOLONTAIRES) FURENT ALÉATOIREMENT RÉPARTIS EN **DEUX GROUPES**.



LE GROUPE 1 PRENAIT UN **PLACEBO**, UNE PILULE IDENTIQUE À L'ASPIRINE MAIS SANS ASPIRINE.



LE GROUPE 2 RECEVAIT UNE ASPIRINE PAR JOUR.

SUR UNE PÉRIODE AVOISINANT LES CINQ ANS*,
LES ENQUÊTEURS ENREGISTRÈRENT LES RÉPONSES :
CRISE CARDIAQUE OU NON.
LE RÉSULTAT (DANS LES NOMBRES
QUI SUIVENT NOUS AVONS SOMMÉ
LES CRISES CARDIAQUES FATALES OU NON) :

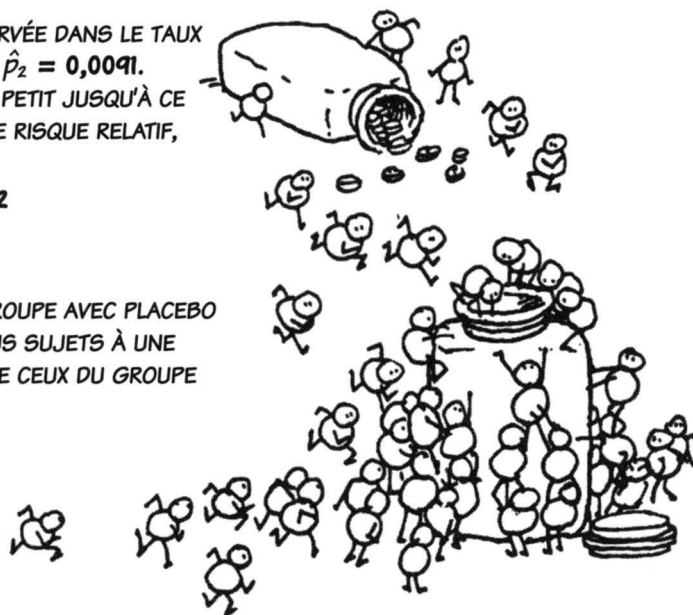


	ATTAQUES	PAS D'ATTAQUES	n	TAUX D'ATTAQUES
PLACEBO	239	10 795	11 034	$\hat{p}_1 = \frac{239}{11\,034} = 0,0217$
ASPIRINE	139	10 898	11 037	$\hat{p}_2 = \frac{139}{11\,037} = 0,0126$

LA DIFFÉRENCE OBSERVÉE DANS LE TAUX
DE SUCCÈS EST $\hat{p}_1 - \hat{p}_2 = 0,0091$.
CELA PEUT PARAÎTRE PETIT JUSQU'À CE
QUE L'ON REGARDE LE RISQUE RELATIF,

$$\frac{\hat{p}_1}{\hat{p}_2} = \frac{0,0217}{0,0126} = 1,72$$

LES MEMBRES DU GROUPE AVEC PLACEBO
ÉTAIENT 1,72 FOIS PLUS SUJETS À UNE
CRISE CARDIAQUE QUE CEUX DU GROUPE
AVEC L'ASPIRINE.



*L'ÉTUDE FUT ARRÊTÉE AVANT SON TERME À CAUSE DE SON RÉSULTAT POSITIF. IL AURAIT ÉTÉ IMPRUDENT
ET IRRÉALISTE DE CACHER LES RÉSULTATS AU GROUPE QUI PRENAIT LE PLACEBO.

Le modèle : LES OBSERVATIONS DES GROUPES « PLACEBO ET ASPIRINE » SONT DES ÉCHANTILLONS INDÉPENDANTS DE DEUX POPULATIONS BINOMIALES. PAR SOUCI DE COHÉRENCE, NOUS POSONS QU'UNE CRISE CARDIAQUE DÉFINIT UN **SUCCÈS** (!).



POPULATION 1 PLACEBO
PROBABILITÉ DE SUCCÈS = p_1



POPULATION 2 ASPIRINE
PROBABILITÉ DE SUCCÈS = p_2

L'OBJECTIF EST D'ESTIMER LA VRAIE DIFFÉRENCE $p_1 - p_2$.

POUR CHAQUE POPULATION (EN FAIT POUR DE GRANDS ÉCHANTILLONS DE LA POPULATION GÉNÉRALE), NOUS AVONS LES VARIABLES ALÉATOIRES USUELLES :

X_1 NOMBRE DE SUCCÈS
DE L'ÉCHANTILLON 1

X_2 NOMBRE DE SUCCÈS
DE L'ÉCHANTILLON 2

$\hat{p}_1 = \frac{X_1}{n_1}$ PROPORTION DE SUCCÈS
DE L'ÉCHANTILLON 1

$\hat{p}_2 = \frac{X_2}{n_2}$ PROPORTION DE SUCCÈS
DE L'ÉCHANTILLON 2

ET UN ESTIMATEUR DU TAUX DE DIFFÉRENCE QUI EST $\hat{p}_1 - \hat{p}_2$.

ET MAINTENANT, TEL UN DISQUE RAYÉ
NOUS NOUS DEMANDONS : COMMENT
 $\hat{p}_1 - \hat{p}_2$ EST-IL DISTRIBUÉ ?



COMMENT ?

COMMENT ?

COMMENT ?

Distribution d'échantillonnage de $\hat{P}_1 - \hat{P}_2$

POUR DE GRANDS ÉCHANTILLONS, LA DISTRIBUTION DE $\hat{P}_1 - \hat{P}_2$ EST APPROXIMATIVEMENT NORMALE COMME DANS LE CAS D'UN UNIQUE ÉCHANTILLON. ON PEUT TRANSFORMER EN Z POUR OBTENIR (APPROXIMATIVEMENT) UNE LOI NORMALE STANDARD

$$z = \frac{\hat{P}_1 - \hat{P}_2 - (p_1 - p_2)}{\sigma(\hat{P}_1 - \hat{P}_2)}$$

MAIS COMMENT TROUVER L'ÉCART-TYPE DU DÉNOMINATEUR ?



COMME LES DEUX ÉCHANTILLONS SONT INDÉPENDANTS, LES VARIABLES ALÉATOIRES \hat{P}_1 ET \hat{P}_2 LE SONT AUSSI ET LEURS VARIANCES S'ADDITIONNENT.

$$\sigma^2(\hat{P}_1 - \hat{P}_2) = \sigma^2(\hat{P}_1) + \sigma^2(\hat{P}_2)$$

DONC

$$\sigma(\hat{P}_1 - \hat{P}_2) = \sqrt{\sigma^2(\hat{P}_1) + \sigma^2(\hat{P}_2)}$$

MAINTENANT QUE NOUS CONNAISSONS LA DISTRIBUTION DE LA STATISTIQUE DE TEST, ON PEUT ESTIMER DES **INTERVALLES DE CONFIANCE** ET **TESTER L'HYPOTHÈSE** QUE L'ASPIRINE RÉDUIT LE RISQUE DE CRISE CARDIAQUE.



Intervalle de confiance de $p_1 - p_2$

COMME D'HABITUDE, L'INTERVALLE
DE CONFIANCE POUR NOTRE ESTIMATION
EST DU TYPE :

$$p_1 - p_2 \in \hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} s(\hat{p}_1 - \hat{p}_2)$$

\nearrow VRAIE DIFFÉRENCE DES PROPORTIONS DE POPULATION
 \uparrow DIFFÉRENCE OBSERVÉE
 \uparrow FRACTILE OU VALEUR CRITIQUE
 \nwarrow ERREUR-TYPE ESTIMÉE

LES VARIANCES DE \hat{p}_1 ET \hat{p}_2 S'AJOUTENT, DONC
L'ERREUR-TYPE DEVIENT :

$$s(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$$

DANS L'ÉTUDE SUR L'ASPIRINE, L'ERREUR-TYPE EST :

$$\sqrt{\frac{(0,0217)(0,9783)}{11\,034} + \frac{(0,0126)(0,9874)}{11\,037}}$$

$$= 0,00175$$



DANS LE CAS DE L'ASPIRINE,
POUR OBTENIR UN INTERVALLE
DE CONFIANCE À 95 %, ON REMPLACE
PAR LES VALEURS :

$$p_1 - p_2 \in 0,0091 \pm (1,96)(0,00175)$$

$$p_1 - p_2 \in 0,0091 \pm 0,0034$$



TRADUCTION :

NOUS SOMMES CONFIANTS
AU MOINS À 95 % QUE LA DIFFÉRENCE
DANS LE TAUX DE CRISES CARDIAQUES
EST ENTRE **0,0057** ET **0,0125**,
UN NOMBRE DÉFINITIVEMENT POSITIF!
NOUS SOMMES DONC SÛRS
AU MOINS À 95 % QUE L'ASPIRINE
ABAISSSE VRAIMENT LE TAUX DE CRISES
CARDIAQUES.

HUM!
VOUDRIEZ-VOUS
AJOUTER DE L'ASPIRINE
À MES CROQUETTES?



Tests d'hypothèses

LA QUESTION FORMELLE D'UN TEST D'HYPOTHÈSE EST :



H_0 : L'HYPOTHÈSE NULLE EST QUE L'ASPIRINE N'A PAS D'EFFET : $p_1 = p_2$.

H_a : L'ALTERNATIVE EST QUE L'ASPIRINE RÉDUIT LE TAUX DE CRISES CARDIAQUES : $p_1 > p_2$.

MAINTENANT, IL NOUS FAUT UNE STATISTIQUE DISTRIBUÉE NORMALEMENT POUR LAQUELLE H_0 EST VRAIE.



NOTEZ QUE SOUS H_0 , LES DEUX PROPORTIONS SONT ÉGALES $p_1 = p_2 = p \dots$ ON MÉLANGE LES DONNÉES POUR AVOIR LA PROPORTION DE CRISES CARDIAQUES DANS LES DEUX ÉCHANTILLONS PRIS ENSEMBLE :

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$$

QUAND L'HYPOTHÈSE NULLE EST VRAIE, L'ERREUR-TYPE NOTÉE s_0 (INDICE 0 POUR H_0 VRAIE) DÉPEND UNIQUEMENT DE L'ESTIMATION MUTUALISÉE :

$$s_0(\hat{p}_1 - \hat{p}_2) = \sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

ET ON A CETTE STATISTIQUE DE TEST :

$$z = \frac{\hat{p}_1 - \hat{p}_2}{s_0(\hat{p}_1 - \hat{p}_2)}$$

(LE NUMÉRATEUR EST EN GÉNÉRAL PLUTÔT $\hat{p}_1 - \hat{p}_2 - (p_1 - p_2)$, MAIS H_0 SUPPOSE QUE $p_1 - p_2 = 0$.)



POUR L'ÉTUDE DES ASPIRINES, ON TROUVE :

$$\hat{p} = \frac{378}{22071} = 0,0171 \text{ ET } 1 - \hat{p} = 0,9829$$

$$s_0(\hat{p}_1 - \hat{p}_2) = 0,00175 \text{ DONC :}$$

$$z_{\text{OBS}} = \frac{\hat{p}_1 - \hat{p}_2}{s_0(\hat{p}_1 - \hat{p}_2)} = \frac{0,0041}{0,00175}$$

$$= 5,20$$

z_{OBS} EST À PLUS DE **CINQ ÉCARTS-TYPES** DE ZÉRO, L'EFFET EST FORTEMENT POSITIF. EN UTILISANT UNE TABLE OU UN ORDINATEUR, ON OBTIENT LA VALEUR P :

$$\text{VALEUR DE } P = P(z \geq z_{OBS}) = P(z \geq 5,2) = 0,0000001$$



SI L'HYPOTHÈSE NULLE ÉTAIT VRAIE, LA PROBABILITÉ D'OBSERVER UN EFFET AUSSI LARGE EST DE UN SUR **DIX MILLIONS** – UNE PREUVE SOLIDE **CONTRE H_0** !

La recette générale :



POUR TESTER L'HYPOTHÈSE NULLE :

$$H_0: p_1 = p_2$$

ON CALCULE LA STATISTIQUE DE TEST :

$$z_{OBS} = \frac{\hat{p}_1 - \hat{p}_2}{s_0(\hat{p}_1 - \hat{p}_2)}$$

OÙ s_0 EST CALCULÉE EN UTILISANT LA PROBABILITÉ MUTUALISÉE EN MÉLANGEANT LES GROUPES.



LA VALEUR P PERTINENTE DÉPEND DE L'HYPOTHÈSE ALTERNATIVE :

A) **BILATÉRAL**, $H_a: p_1 \neq p_2$



$$\text{VALEUR DE } P = P(|z| > z_{OBS})$$

B) **UNILATÉRAL DROIT**, $H_a: p_1 > p_2$



$$\text{VALEUR DE } P = P(z > z_{OBS})$$

C) **UNILATÉRAL GAUCHE**, $H_a: p_1 < p_2$



$$\text{VALEUR DE } P = P(z < z_{OBS})$$

L'ANALYSE DE L'ÉTUDE SUR L'ASPIRINE DÉPENDAIT DE CERTAINES CARACTÉRISTIQUES DE L'EXPÉRIENCE CONÇUES POUR GARANTIR LE HASARD ET ÉLIMINER LES BIAIS.



LES POINTS 1 ET 2 SONT ESSENTIELS DANS LA CONCEPTION DE LA PLUPART DES ESSAIS CLINIQUES. MAIS LE POINT 3 N'EST PAS INDISPENSABLE. IL EXISTE DE BONS TESTS SUR DE PETITS ÉCHANTILLONS, DISPONIBLES DANS DES LOGICIELS DE STATISTIQUES. CES PROCÉDURES **NON PARAMÉTRIQUES** DÉPENDENT DE CALCULS DE PROBABILITÉS SIMPLES MAIS LONGS DU TYPE DE CEUX QUE NOUS AVONS RENCONTRÉS DANS LE CHAPITRE 4.

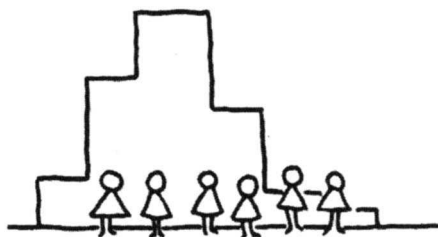


Comparaison de MOYENNES de deux populations

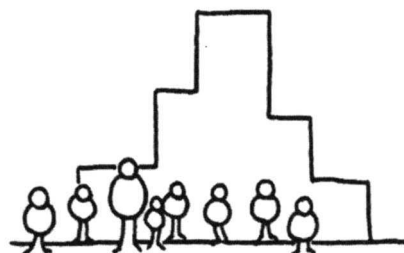
SUPPOSONS QUE L'ON VEUILLE COMPARER
LE SALAIRE MOYEN DES HOMMES À CELUI
DES FEMMES EMPLOYÉES POUR LE MÊME
TRAVAIL DANS UNE ENTREPRISE.



LA POPULATION 1 EST CELLE DES FEMMES, ET LA POPULATION 2 CELLE DES HOMMES.



LA POPULATION 1 A UN SALAIRE
MOYEN μ_1 ET UN ÉCART-TYPE σ_1 .



LA POPULATION 2 A UN SALAIRE
MOYEN μ_2 ET UN ÉCART-TYPE σ_2 .

UN ÉCHANTILLON ALÉATOIRE SIMPLE DE TAILLE n_1 POUR LE GROUPE 1 ET DE TAILLE n_2 POUR LE GROUPE 2 FOURNIT DES MOYENNES D'ÉCHANTILLONS DE \bar{x}_1 ET \bar{x}_2 AVEC DES ÉCARTS-TYPES RESPECTIVEMENT DE s_1 ET s_2 . L'ESTIMATEUR DE $\mu_1 - \mu_2$ EST :

$$\bar{X}_1 - \bar{X}_2$$

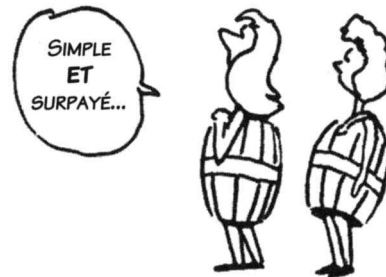
À QUEL POINT NOTRE ESTIMATEUR $\bar{X}_1 - \bar{X}_2$
EST-IL BON ?

POUR DE GRANDS ÉCHANTILLONS,
C'EST PRESQUE NORMAL D'APRÈS
LE THÉORÈME CENTRAL LIMITE
ET L'ERREUR-TYPE VAUT

$$s(\bar{X}_1 - \bar{X}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

(LES VARIANCES S'AJOUTENT
CAR LES ÉCHANTILLONS SONT INDÉPENDANTS.)
MAINTENANT NOUS POUVONS DIRECTEMENT
POURSUIVRE AVEC LES **intervalles**
de confiance : POUR DE GRANDS
ÉCHANTILLONS, L'INTERVALLE AVEC $(1 - \alpha)$
DE CONFIANCE POUR LA DIFFÉRENCE
DE MOYENNE EST :

$$\mu_1 - \mu_2 \in \bar{x}_1 - \bar{x}_2 \pm z_{\alpha/2} s(\bar{X}_1 - \bar{X}_2)$$



Test d'hypothèse : ON POSE L'HYPOTHÈSE NULLE QUI DIT QUE LES MOYENNES
DE POPULATION SONT ÉGALES.

$$H_0 : \mu_1 = \mu_2$$

LA STATISTIQUE DE TEST VAUT :

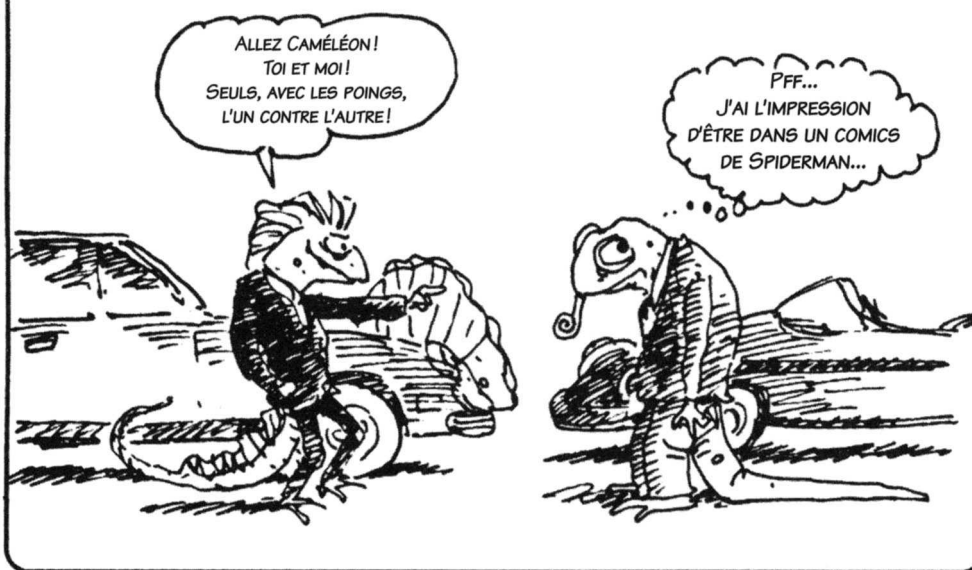
$$z_{OBS} = \frac{\bar{X}_1 - \bar{X}_2}{s(\bar{X}_1 - \bar{X}_2)}$$

ET LES VALEURS P FONCTIONNENT
DE MANIÈRE HABITUELLE.



Et qu'en est-il de la comparaison des MOYENNES DE PETITS ÉCHANTILLONS?

VOUS VOUS RAPPELEZ CAMÉLÉON AUTOMOBILES? LEUR CONCURRENT, IGUANE AUTOS, AFFIRME QUE LEUR CAPOT AVEC DÉCORATION EN POLYSTYRÈNE OFFRE UNE MEILLEURE PROTECTION AUX CRASHS FRONTAUX ET ILS ONT CRASHÉ SEPT IGUANE POUR LE PROUVER!



LEURS RÉSULTATS COMPARÉS AUX CAMÉLÉON SONT :

CAMÉLÉON		IGUANE	
1	150 €	1	50 €
2	400 €	2	200 €
3	720 €	3	150 €
4	500 €	4	400 €
5	930 €	5	750 €
		6	400 €
		7	150 €
n_1	5	n_2	7
\bar{x}_1	540 €	\bar{x}_2	300 €
s_1	299 €	s_2	238 €



LE t DE STUDENT PEUT ÊTRE UTILISÉ SI LES DEUX POPULATIONS ONT UNE FORME DE MONTICULE ET ONT LE MÊME ÉCART-TYPE $\sigma = \sigma_1 = \sigma_2$. LE SEUL HIC EST QUE NOUS DEVONS FAIRE UNE **MOYENNE PONDÉRÉE** DES VARIANCES D'ÉCHANTILLON POUR ESTIMER σ^2 .

$$s_{\text{EST}}^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$



L'ERREUR-TYPE EST LA MÊME QUE POUR DE GRANDS ÉCHANTILLONS, SAUF QUE s_{EST} REMPLACE s_1 ET s_2 .

$$s(\bar{X}_1 - \bar{X}_2) = \sqrt{\frac{s_{\text{EST}}^2}{n_1} + \frac{s_{\text{EST}}^2}{n_2}}$$

$$= s_{\text{EST}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

L'INTERVALLE AVEC $1 - \alpha$ DE CONFIANCE EST ALORS :

$$\mu_1 - \mu_2 \in \bar{x}_1 - \bar{x}_2 \pm t_{\alpha/2} s(\bar{X}_1 - \bar{X}_2)$$

OÙ $t_{\alpha/2}$ EST LE FRACTILE D'UN t DE STUDENT AVEC $n_1 + n_2 - 2$ DEGRÉS DE LIBERTÉ.

LES CONSTRUCTEURS REPTILIENS D'AUTOMOBILE CONVIENNENT QUE LEURS ÉCARTS-TYPES SONT PROCHES ET QUE LEURS HISTOGRAMMES DE RÉPARATIONS ONT UNE FORME DE CLOCHE. ILS CALCULENT :

$$s_{\text{EST}} = \sqrt{\frac{4 \times 299^2 + 6 \times 238^2}{10}} = 264$$

$$s(\bar{X}_1 - \bar{X}_2) = 264 \sqrt{\frac{1}{5} + \frac{1}{7}} = 155$$

L'INTERVALLE DE CONFIANCE À 95 % EST :

$$\mu_1 - \mu_2 \in 540 - 300 \pm t_{0,025} \times 155$$

$$\mu_1 - \mu_2 \in 240 \pm 2,23 \times 155$$

$$\mu_1 - \mu_2 \in 240 \pm 346$$

COMME CELA INCLUT LA VALEUR 0, IGUANE AUTOS N'A PAS PROUVÉ LA PRÉSENCE D'AMÉLIORATIONS SIGNIFICATIVES POUR LES COÛTS DE RÉPARATION.

OK. OUBLIONS LA SÉCURITÉ...
VOUS NE FAITES PAS LE POIDS
AU NIVEAU DU STYLE...



LE PROCHAIN EXEMPLE MONTRE L'INCONVÉNIENT
DE SUIVRE BÊTEMENT LE LIVRE DE RECETTES :
LE PROPRIÉTAIRE D'UNE LARGE FLOTTE DE TAXIS
AMÉRICAINS VEUT COMPARER LES DISTANCES
PARCOURUES, EN MILES, AVEC, SOIT
UNE **ESSENCE A**, SOIT UNE **ESSENCE B**.



IL SÉLECTIONNE 100 TAXIS, ET LEUR ASSIGNE AU HASARD L'UNE DES ESSENCES.
APRÈS UNE JOURNÉE TYPE DE TAXI, IL OBTIENT :

	TAILLE D'ÉCHANTILLON	MOYENNE DE MILES	ÉCART-TYPE
A	50	25	5,00
B	50	26	4,00



LA DIFFÉRENCE D'ÉCHANTILLON EST :

$$\bar{x}_1 - \bar{x}_2 = 25 - 26 = -1$$

L'ESSENCE B EST-ELLE
MEILLEURE QUE L'ESSENCE A ?



EN RAISON DES ÉCARTS-TYPES ÉLEVÉS,
L'ERREUR-TYPE EST ASSEZ SUBSTANTIELLE.

$$s(\bar{X}_1 - \bar{X}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$= \sqrt{\frac{25}{50} + \frac{16}{50}}$$

$$= 0,905$$

AVEC 95 % DE CONFIANCE, NOUS AVONS :

$$\mu_1 - \mu_2 \in \bar{x}_1 - \bar{x}_2 \pm z_{0,025} s(\bar{X}_1 - \bar{X}_2)$$

$$\mu_1 - \mu_2 \in -1 \pm (1,96 \times 0,905)$$

$$\mu_1 - \mu_2 \in -1 \pm 1,774$$

CELA INCLUT LA VALEUR 0, QUI CORRESPOND
À $\mu_1 = \mu_2$.

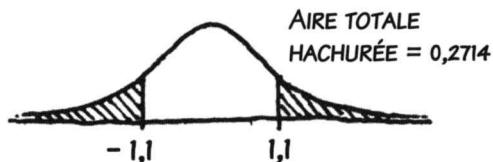


LA VALEUR P POUR L'HYPOTHÈSE
ALTERNATIVE, $H_a: \mu_1 \neq \mu_2$, VAUT :

$$P(|z| > |z_{\text{obs}}|) = P\left(|z| > \frac{1}{0,905}\right)$$

$$= P(|z| > 1,1) = 2 \times 0,1357$$

$$= 0,2714$$



CELA EXCÈDE LE SEUIL DE SIGNIFICATION
 $\alpha = 0,05$. NOUS EN CONCLUONS
DONC QUE LES PREUVES EN FAVEUR
DE L'ESSENCE B SONT TRÈS FAIBLES.



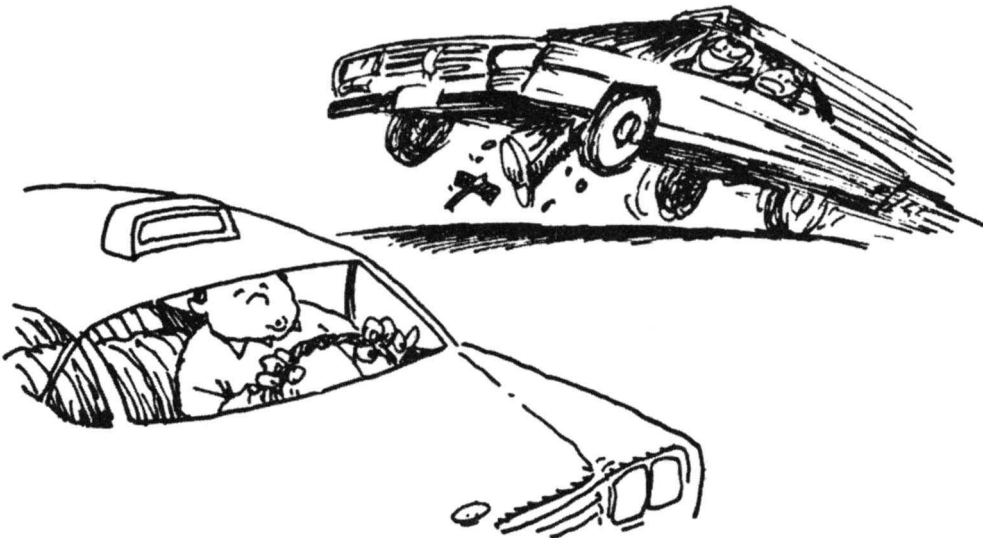
COMPARAISON APPARIÉE

Une meilleure façon pour comparer les essences



LE PROPRIÉTAIRE DE TAXIS A SUIVI LE LIVRE À LA LETTRE. SES ÉCHANTILLONS ÉTAIENT ALÉATOIRES D'UNE TAILLE SUFFISAMMENT GRANDE. IL A SIMPLEMENT OUBLIÉ DE **RÉFLÉCHIR** QUAND C'ÉTAIT NÉCESSAIRE.

BIEN QUE L'ESSENCE B SEMBLE LÉGÈREMENT SUPÉRIEURE À L'ESSENCE A, L'INTERVALLE DE CONFIANCE ÉTAIT LARGE À CAUSE DES GRANDS ÉCARTS-TYPES (AINSI, **LES KILOMÈTRES PARCOURUS VARIAIENT BEAUCOUP D'UN TAXI À L'AUTRE**). POURQUOI TANT DE VARIABILITÉ? PARCE QUE LES TAXIS ET LEURS CONDUCTEURS ONT DES **CARACTÉRISTIQUES** DIFFÉRENTES.



UNE BIEN MEILLEURE FAÇON D'ABORDER CETTE ÉTUDE EST D'AFFECTER L'ESSENCE A ET B AU MÊME TAXI SUR DES JOURS DIFFÉRENTS.



LE TRAITEMENT PEUT ENCORE ÊTRE ALÉATOIRE AVEC UN LANCER DE PIÈCE POUR DÉCIDER SI ON UTILISE L'ESSENCE A LE MARDI OU LE MERCREDI. ET ON PEUT MÊME ÉCONOMISER L'ARGENT ET LE TEMPS DU PROPRIÉTAIRE EN NE PRATIQUANT L'EXPÉRIENCE QUE SUR 10 TAXIS !



TAXI	ESSENCE A	ESSENCE B	DIFFÉRENCE
1	27,01	26,95	0,06
2	20,00	20,44	- 0,44
3	23,41	25,05	- 1,64
4	25,22	26,32	- 1,10
5	30,11	29,56	0,55
6	25,55	26,60	- 1,05
7	22,23	22,93	- 0,70
8	19,78	20,23	- 0,45
9	33,45	33,95	- 0,50
10	25,22	26,01	- 0,79
MOYENNE	25,198	25,804	- 0,61
ÉCART-TYPE	4,27	4,10	0,61

NOTEZ QUE LES MOYENNES ET ÉCARTS-TYPES DES ESSENCES A ET B SONT À PEU PRÈS LES MÊMES QUE PRÉCÉDEMMENT. C'ÉTAIT PRÉVISIBLE PUISQUE CES STATISTIQUES ONT LES MÊMES SOURCES DE VARIATION QUE DES ÉCHANTILLONS NON APPARIÉS. MAIS MAINTENANT LA **COLONNE DES DIFFÉRENCES** A UN FAIBLE ÉCART-TYPE. CES DIFFÉRENCES ÉLIMINENT LA **VARIABILITÉ ENTRE LES TAXIS** EN COMPARANT LES PERFORMANCES DE L'ESSENCE POUR UNE MÊME VOITURE.

LA DIFFÉRENCE d_i FOURNIT UNE MESURE DE LA DIFFÉRENCE POUR CHAQUE TAXI ET NOUS POUVONS L'UTILISER POUR NOTRE STATISTIQUE DE TEST t SUR DE PETITS ÉCHANTILLONS.

$$t = \frac{\bar{d}}{s_d/\sqrt{n}}$$

OÙ \bar{d} EST LA MOYENNE DES DIFFÉRENCES DE L'ÉCHANTILLON (- 0,61 DANS L'EXEMPLE) ET s_d EST LEUR ÉCART-TYPE (ICI 0,61).

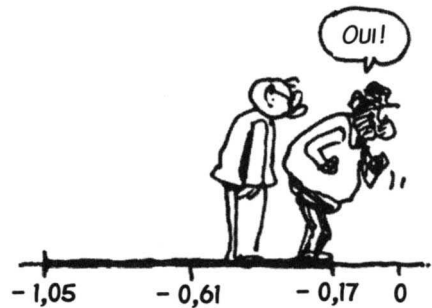


L'INTERVALLE DE CONFIANCE À 95 % CENTRÉ EN \bar{d} EST :

$$\mu_d \in \bar{d} \pm t_{0,025} \times \frac{s_d}{\sqrt{n}}$$

$$\mu_d \in -0,61 \pm 2,26 \times \frac{0,61}{\sqrt{10}}$$

$$\mu_d \in -0,61 \pm 0,44$$



NOUS SOMMES SÛRS À 95 % QUE $-1,05 \leq \mu_d \leq -0,17$, CE QUI CONSTITUE UNE PREUVE SOLIDE QUE L'ESSENCE B EST SUPÉRIEURE.

LA VALEUR P DU TEST D'HYPOTHÈSE PEUT ÊTRE CALCULÉE PAR ORDINATEUR.

$$H_a: \mu_d \neq 0$$

$$\text{VALEUR DE } P = P(|t| \geq |t_{\text{obs}}|)$$

$$= P\left(|t| \geq \frac{0,61}{0,19}\right)$$

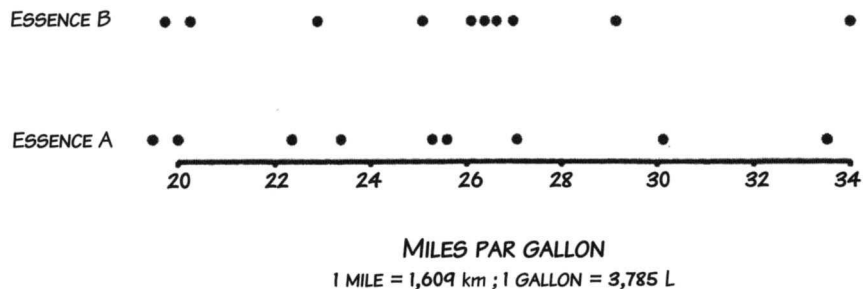
$$= P(|t| \geq 3,21)$$

$$= 0,011$$

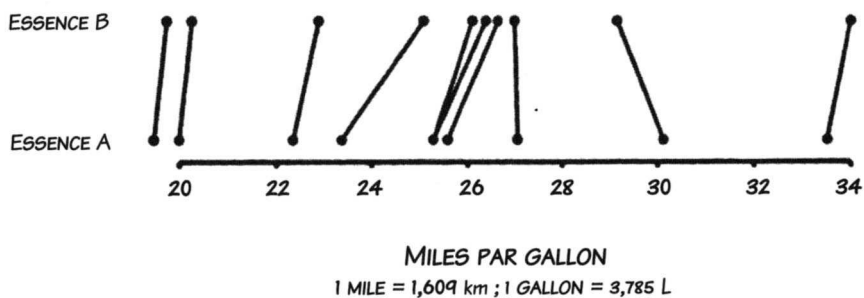


LA VALEUR P EST INFÉRIEURE À 0,05, DONC, À NOUVEAU, L'ESSENCE B PASSE LE TEST.

VOICI LE DIAGRAMME EN POINTS DES DONNÉES EN MILES PARCOURUS SELON L'ESSENCE : LE PREMIER DESSIN MONTRE LES DONNÉES NON APPARIÉES.



ET VOICI LES MÊMES DONNÉES APPARIÉES PAR TAXI.



LA PRÉDOMINANCE
DE SEGMENTS INCLINÉS VERS
LA DROITE EST UN INDICE FORT
QUE L'ESSENCE B ACCROÎT
LA DISTANCE PARCOURUE.



UTILISER DES ÉCHANTILLONS APPARIÉS EST LA MÉTHODE LA PLUS EFFICACE POUR RÉDUIRE LA VARIABILITÉ NATURELLE QUAND ON COMPARE DES TRAITEMENTS. PAR EXEMPLE, SI L'ON COMPARE DEUX MARQUES DE CRÈMES POUR LES MAINS, ON PEUT ALÉATOIREMENT AFFECTER LA MAIN DROITE À L'UNE ET LA MAIN GAUCHE À L'AUTRE POUR LE MÊME SUJET. ON ÉLIMINE AINSI LA VARIABILITÉ DUE À DES TYPES DE PEaux DIFFÉRENTS.



OU SI ON COMPARE DEUX MARQUES DE CÉRÉALES, CHAQUE « GOÛTEUR » ÉVALUERA LES DEUX CÉRÉALES (DANS UN ORDRE ALÉATOIRE). UNE COMPARAISON APPARIÉE SUPPRIME LE BIAIS NATUREL D'UN GOÛTEUR QUI SERAIT POUR OU CONTRE LES CÉRÉALES D'UNE MANIÈRE GÉNÉRALE.



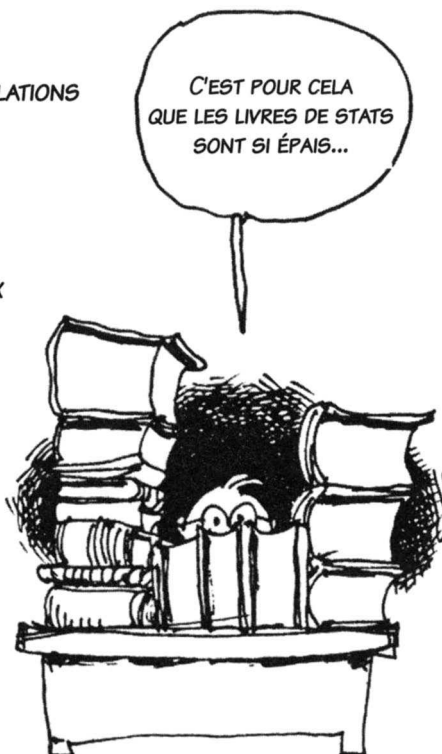
DANS CE CHAPITRE, NOUS AVONS APPLIQUÉ LES IDÉES DE BASE SUR LES INTERVALLES DE CONFIANCE ET LES TESTS D'HYPOTHÈSES À LA COMPARAISON DE DEUX POPULATIONS. IL EXISTE D'INNOMBRABLES AUTRES POSSIBILITÉS. NOUS AURIONS PU POURSUIVRE EN DÉCRIVANT DES COMPARAISONS :

- D'ÉCARTS-TYPES DE DEUX POPULATIONS AVEC DE PETITS ÉCHANTILLONS;

- DE MOYENNES DE PLUS DE DEUX POPULATIONS AVEC DE GRANDS ÉCHANTILLONS;

- DE MOYENNES DE PLUS DE DEUX POPULATIONS AVEC DE PETITS ÉCHANTILLONS;

ETC.!



DANS LA PRATIQUE, LES STATISTICIENS DÉTERMINENT LA NATURE GÉNÉRALE D'UN PROBLÈME ET CONSULTENT ALORS LE LIVRE DE RÉFÉRENCE ADAPTÉ.



LA SEULE IDÉE VÉRITABLEMENT NOUVELLE DE CE CHAPITRE ÉTAIT L'**APPARIEMENT** POUR RÉALISER DES TESTS DE **COMPARAISON**. DANS LE PROCHAIN CHAPITRE, NOUS VERRONS D'AUTRES TYPES DE DISPOSITIFS EXPÉRIMENTAUX.

VOUS VOULEZ ACHETER
UN CAMÉLÉON
D'OCCASION ?



Chapitre 10

Méthodes expérimentales

LA CONCEPTION EST SOUVENT CE QUI GARANTIT LE SUCCÈS OU L'ÉCHEC D'UNE EXPÉRIENCE. DANS L'EXEMPLE DE LA COMPARAISON APPARIÉE, NOTRE STATISTICIEN A INVERSÉ LES RÔLES PASSANT D'UNE COLLECTE ET D'UNE ANALYSE PASSIVE DES DONNÉES À UNE PARTICIPATION ACTIVE DANS LE DISPOSITIF DE L'EXPÉRIENCE.

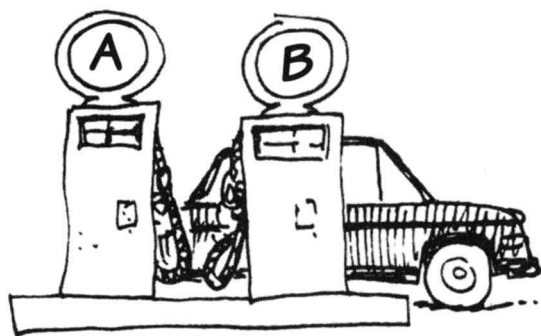
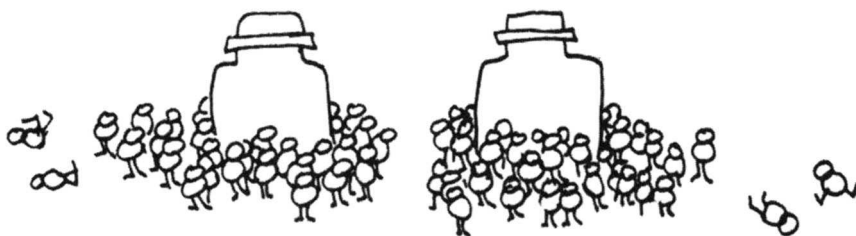


DANS CE CHAPITRE,
NOUS INTRODUISONS
LES IDÉES DE BASE
DES MÉTHODES
EXPÉRIMENTALES,
EN LAISSANT LES ANALYSES
NUMÉRIQUES DÉTAILLÉES
À VOTRE LOGICIEL
DE STATISTIQUES.



PAS DE FORMULES
DANS CE CHAPITRE...
DÉSOLÉ!

LES ÉLÉMENTS D'UN DISPOSITIF SONT : LES **UNITÉS EXPÉRIMENTALES**
ET LES TRAITEMENTS QUI SERONT ASSIGNÉS AUX UNITÉS. L'OBJECTIF DU DISPOSITIF
EST DE COMPARER LES TRAITEMENTS.



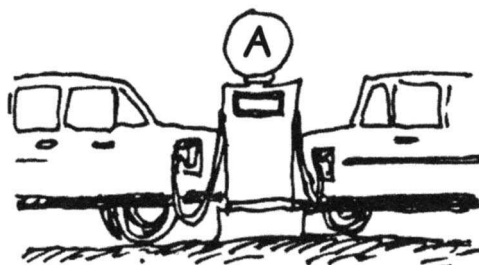
POUR DES ESSAIS
MÉDICAUX, LES **PATIENTS**
SONT LES UNITÉS,
ET LES **MÉDICAMENTS**
SONT LES TRAITEMENTS.
DANS L'EXEMPLE SUR
L'ESSENCE, LES UNITÉS
EXPÉRIMENTALES
SONT LES TAXIS,
ET LES TRAITEMENTS
À COMPARER SONT
LES ESSENCES A ET B.

DANS DES EXPÉRIENCES AGRICOLES, LES UNITÉS EXPÉRIMENTALES SONT SOUVENT
DES PARCELLES DE TERRAIN, ET LES TRAITEMENTS PEUVENT CORRESPONDRE
À L'APPLICATION DE DIFFÉRENTES VARIÉTÉS DE BLÉ, DE PESTICIDES, DE FERTILISANTS, ETC.

AUJOURD'HUI, LES IDÉES DES MÉTHODES EXPÉRIMENTALES SONT LARGEMENT UTILISÉES DANS L'**OPTIMISATION DES PROCESSUS INDUSTRIELS**, EN **MÉDECINE** ET DANS LES **SCIENCES SOCIALES**. TOUT DISPOSITIF EXPÉRIMENTAL UTILISE **TROIS PRINCIPES DE BASE** QUI ÉTAIENT TOUS PRÉSENTS DANS L'EXEMPLE DES TAXIS.



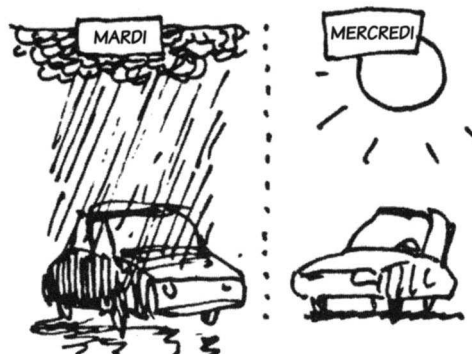
La réplication : LES MÊMES TRAITEMENTS SONT ASSIGNÉS À DIFFÉRENTES UNITÉS EXPÉRIMENTALES. SANS RÉPLICATION, IL EST IMPOSSIBLE D'ÉVALUER LA VARIABILITÉ NATURELLE ET LES ERREURS DE MESURE.



Le contrôle local : IL SE RÉFÈRE À TOUTE MÉTHODE QUI TIEN COMTE DE LA VARIABILITÉ NATURELLE ET LA RÉDUIT. UN MOYEN EST DE REGROUPER LES UNITÉS EXPÉRIMENTALES SIMILAIRES EN BLOCS. DANS L'EXEMPLE DES TAXIS, CHACUN D'EUX UTILISAIT LES DEUX ESSENCES, ON DIT ALORS QUE LE TAXI EST UN BLOC.



La randomisation : C'EST L'ÉTAPE ESSENTIELLE POUR TOUTE ÉTUDE STATISTIQUE! LES TRAITEMENTS DOIVENT ÊTRE ASSIGNÉS ALÉATOIREMENT AUX UNITÉS EXPÉRIMENTALES. POUR CHAQUE TAXI, L'ESSENCE A EST UTILISÉE LE MARDI OU LE MERCREDI À PILE OU FACE. SI CE N'AVAIT PAS ÉTÉ LE CAS, LES RÉSULTATS AURAIENT PU ÊTRE FAUSSÉS PAR DES DIFFÉRENCES ENTRE LE MARDI ET LE MERCREDI.



SUPPOSONS MAINTENANT QUE L'ON VEUILLE ÉTUDIER LES EFFETS DE DEUX MARQUES DE PNEUS ET AUSSI DE DEUX ESSENCES. IL Y A QUATRE TRAITEMENTS POSSIBLES QUE NOUS POUVONS PRÉSENTER DANS UNE TABLE FACTORIELLE 2×2 . LES DEUX FACTEURS ÉTANT L'ESSENCE ET LE TYPE DE PNEU.

	ESSENCE A	ESSENCE B
PNEU A	a	b
PNEU B	c	d



ON PEUT ASSIGNER LES QUATRE TRAITEMENTS DE FAÇON ALÉATOIRE SUR QUATRE JOURS DIFFÉRENTS POUR CHAQUE TAXI. LES QUATRE TRAITEMENTS (a, b, c et d) SONT RÉPÉTÉS DANS CHAQUE BLOC (TAXI). IL S'AGIT ALORS D'UN **DISPOSITIF EN BLOCS ALÉATOIRES COMPLETS**.

JUSQU'À PRÉSENT, NOUS AVONS FAIT L'HYPOTHÈSE QUE CHAQUE JOUR DE LA SEMAINE EST IDENTIQUE. MAIS ON PEUT CONTRÔLER CELA AUSSI EN UTILISANT SEULEMENT QUATRE TAXIS ET EN ASSIGNANT UN TRAITEMENT SELON LE PLAN DU TABLEAU DE DROITE.

		JOUR			
		1	2	3	4
TAXI	1	a	b	c	d
	2	b	c	d	a
	3	c	d	a	b
	4	d	a	b	c

REMARQUE :
CHAQUE TRAITEMENT
APPARAÎT UNE SEULE
FOIS DANS CHAQUE
COLONNE!

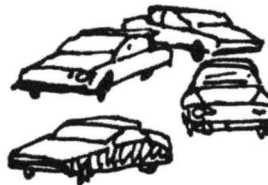


UN TABLEAU 4×4 AVEC 4 ÉLÉMENTS DIFFÉRENTS, CHACUN APPARAISSANT UNE SEULE FOIS DANS CHAQUE LIGNE ET COLONNE, S'APPELLE UN **carré latin**. DANS CETTE EXPÉRIENCE, LES QUATRE JOURS ET LES QUATRE TAXIS REÇOIVENT LES QUATRE TRAITEMENTS EXACTEMENT UNE FOIS CHACUN.



L'ÉTAPE DE RANDOMISATION SÉLECTIONNE UN CARRÉ LATIN AU HASARD DE LA LISTE DE TOUS LES CARRÉS LATINS POSSIBLES À QUATRE ÉLÉMENTS.

SI QUATRE UNITÉS NE SONT PAS SUFFISANTES, ON PEUT AUGMENTER LE NOMBRE D'UNITÉS EXPÉRIMENTALES EN **RÉPÉTANT** LE DISPOSITIF EXPÉRIMENTAL. SI ON PART DE HUIT TAXIS, ON PEUT LES DIVISER EN DEUX GROUPES DE QUATRE ET RÉPÉTER LE DISPOSITIF DANS CHAQUE GROUPE.



NOUS AVIONS PROMIS DE NE PAS ENTRER DANS LE DÉTAIL DE L'ANALYSE DE DONNÉES. MAIS VOICI EN GROS COMMENT GÉRER UN DISPOSITIF COMPLEXE DE CE TYPE.



L'ANALYSE DES DISPOSITIFS EXPÉRIMENTAUX SE FAIT EN ALLOUANT LA VARIABILITÉ TOTALE AUX DIFFÉRENTES SOURCES. DANS L'EXEMPLE DES TAXIS, LES SOURCES DE VARIATION SONT LE TAXI, LA MARQUE DE PNEU, LE TYPE D'ESSENCE, LE JOUR ET UNE ERREUR ALÉATOIRE. L'ANALYSE DE VARIANCE, OU **ANOVA** EN ABRÉGÉ, DIVISE LA VARIATION TOTALE EN PORTIONS POUR CHAQUE SOURCE.

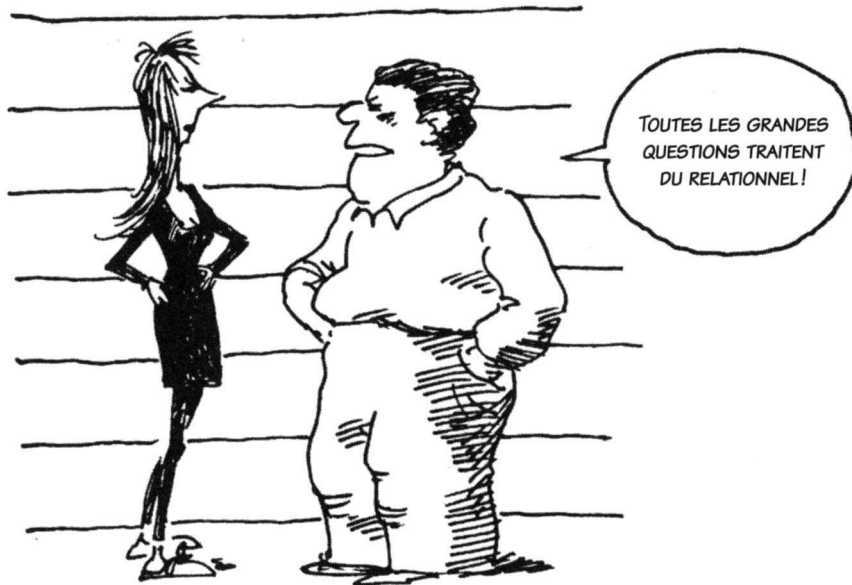
DANS LE PROCHAIN CHAPITRE, NOUS EXPLIQUERONS EN DÉTAIL UN MODÈLE POUR ANALYSER DES DISPOSITIFS COMPLEXES : LE **MODÈLE DE RÉGRESSION LINÉAIRE**. AVEC LA RÉGRESSION LINÉAIRE, VOUS POURREZ VOIR ANOVA NUMÉRIQUEMENT ET DE PRÈS...



Chapitre II

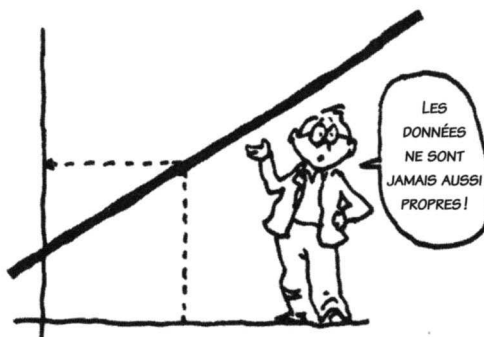
Régression linéaire

JUSQU'À PRÉSENT, NOUS AVONS FAIT DES STATISTIQUES SUR **UNE SEULE VARIABLE À LA FOIS**, QU'ELLE VIENNE D'UNE POPULATION DE PRENEURS DE PILULES OU DE VOITURES ACCIDENTÉES. DANS CE CHAPITRE, NOUS ALLONS VOIR COMMENT RELIER **DEUX** VARIABLES. CONNAISSANT LES POIDS DES 92 ÉTUDIANTS DU CHAPITRE 2, NOUS NOUS DEMANDONS COMMENT ILS SONT RELIÉS À LA TAILLE DE CES ÉTUDIANTS.

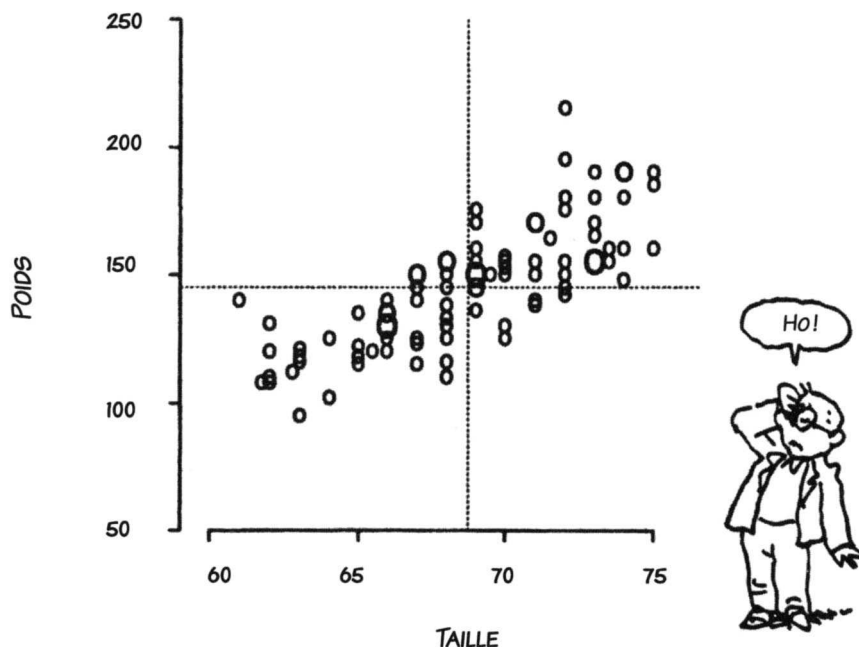


C'EST UN EXEMPLE D'UNE LARGE CLASSE DE QUESTIONS IMPORTANTES : LA **PRESSIION ARTÉRIELLE** PRÉDIT-ELLE LA **DURÉE DE VIE** ? LES **RÉSULTATS DU BAC** PRÉDISSENT-ILS LES **PERFORMANCES EN LICENCE** ? LIRE UN LIVRE DE STATISTIQUES FAIT-IL DE VOUS UNE **MEILLEURE PERSONNE** ?

EN CLASSE DE MATHÉMATIQUES,
VOUS AVEZ PROBABLEMENT VU
DES RELATIONS DÉFINIES COMME
DES GRAPHIQUES. ÉTANT DONNÉ x ,
ON PEUT ALORS PRÉDIRE y .
MAIS EN STATISTIQUES, LES CHOSES
NE SONT PAS AUSSI SIMPLES!
ON SAIT (OU ON SUPPOSE) QUE
LA TAILLE A UNE INCIDENCE SUR
LE POIDS – MAIS CE N'EST PAS
LA SEULE. IL Y A AUSSI D'AUTRES
FACTEURS COMME LE SEXE,
L'ÂGE, LE TYPE PHYSIQUE,
ET UNE **COMPOSANTE ALÉATOIRE**.



POUR CE CHAPITRE, NOTONS y LA DONNÉE DE POIDS EN LIVRES ET x LA DONNÉE DE TAILLE EN POUCES. AINSI (x_i, y_i) EST LA TAILLE ET LE POIDS DE L'ÉTUDIANT i . ON PEUT TRACER LES POINTS (x_i, y_i) DANS UN PLAN, ON APPELLE CELA **UN NUAGE DE POINTS**.



(CERTAINS POINTS SONT PLUS GROS CAR ILS REPRÉSENTENT DEUX OU TROIS ÉTUDIANTS AYANT LE MÊME POIDS ET LA MÊME TAILLE.)

POUVONS-NOUS PRÉVOIR LE POIDS y D'UN ÉTUDIANT CONNAISSANT SA TAILLE x ?

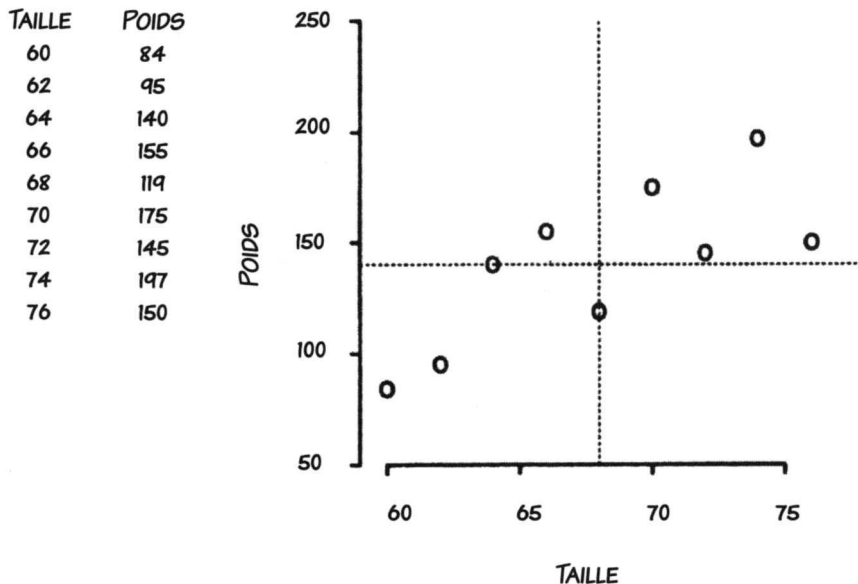
L'ANALYSE DE RÉGRESSION

L'ANALYSE DE RÉGRESSION CONSISTE À AJUSTER UNE LIGNE DROITE À CE NUAGE DÉSORDONNÉ DE POINTS. x EST APPELÉE VARIABLE **INDÉPENDANTE** OU **EXPLICATIVE**, ET y EST LA VARIABLE **DÉPENDANTE** OU DE **RÉPONSE**. LA DROITE DE **RÉGRESSION AFFINE** A LA FORME SUIVANTE :

$$y = ax + b$$



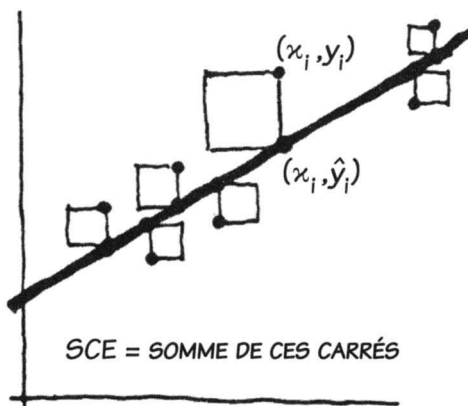
POUR ILLUSTRER LA MÉTHODE D'AJUSTEMENT, UTILISONS UN SOUS-ENSEMBLE UN PEU ARRANGÉ DE DONNÉES CONTENANT SEULEMENT NEUF OBSERVATIONS D'ÉTUDIANTS.



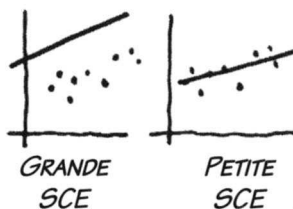
COMMENT OBTENIR LA DROITE LA MIEUX AJUSTÉE AUX DONNÉES ?

L'IDÉE EST DE **MINIMISER** L'ÉCART TOTAL DES VALEURS y OBSERVÉES PAR RAPPORT À LA DROITE. COMME AVEC LA VARIANCE, ON REGARDE PLUTÔT LA **DISTANCE AU CARRÉ** ENTRE y ET LA DROITE, ON LES SOMME POUR OBTENIR LA **SOMME DES CARRÉS DES ERREURS** OU **RÉSIDUS**.

$$SCE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



C'EST UNE MESURE AGRÉGÉE QUI CALCULE DE COMBIEN LA DROITE « PRÉDISANT y_i » (C'EST-À-DIRE \hat{y}_i) DIFFÈRE DES DONNÉES (C'EST-À-DIRE DES VALEURS OBSERVÉES y_i).



LA RÉGRESSION OU DROITE DES MOINDRES CARRÉS

EST LA DROITE AVEC LA PLUS PETITE SCE.



EST-CE QU'ON DOIT LA MESURER POUR CHAQUE DROITE ?

NOTE HISTORIQUE : POURQUOI APPELLE-T-ON CETTE PROCÉDURE L'ANALYSE DE RÉGRESSION ? DANS LES ANNÉES 1880, LE GÉNÉTICIEN FRANCIS GALTON (1822-1911) DÉCOUVRIIT LE PHÉNOMÈNE DE **RÉGRESSION VERS LA MOYENNE**. EN RECHERCHANT LES RÈGLES DE L'HÉRÉDITÉ, IL CONSTATA QUE LA TAILLE DES FILS TENDAIT À **RÉGRESSER** VERS LA TAILLE MOYENNE DE LA POPULATION, PAR RAPPORT À LA TAILLE DES PÈRES. LES PÈRES DE GRANDE TAILLE TENDAIENT À AVOIR DES FILS PLUS PETITS ET VICE VERSA. GALTON DÉVELOPPA L'ANALYSE DE RÉGRESSION POUR ÉTUDIER CE QU'IL APPELA UNE « RÉGRESSION VERS LA MÉDIOCRITÉ ».

GRANDIS MON FILS !



SANS TOURNER AUTOUR DU POT,
NOUS DONNONS LA FORMULE DE LA DROITE
DE RÉGRESSION SANS LA DÉMONSTRER.
C'EST COMPLIQUÉ MAIS CALCULABLE.

$$y = ax + b$$

OÙ

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

ET

$$a = \bar{y} - b\bar{x}$$

OÙ \bar{x} ET \bar{y} SONT LES MOYENNES DE $\{x_i\}$
ET $\{y_i\}$ RESPECTIVEMENT. LES SOMMES
SONT POUR i ALLANT DE 1 À n .



PARCE QUE NOUS ALLONS RETROUVER CERTAINES EXPRESSIONS, ON VA LES ABRÉGER
(DE NOUVEAU, LA SOMME SE FAIT POUR TOUT i , SAUF LORSQUE C'EST PRÉCISÉ).

$$SC_{xx} = \sum (x_i - \bar{x})^2$$

$$SCT = SC_{yy} = \sum (y_i - \bar{y})^2$$

$$SC_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

}

SOMME DES CARRÉS DES ÉCARTS
À LA MOYENNE POUR x_i ET y_i .
SCT VEUT DIRE SOMME DES CARRÉS
TOTAUX.

LE PRODUIT CROISÉ ET SC_{xx}
DÉTERMINENT LE COEFFICIENT b .



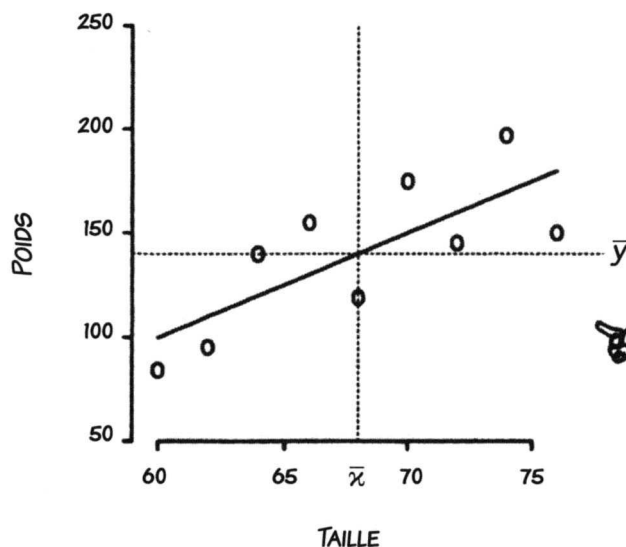
POUR LES DONNÉES ARRANGÉES, VOICI LES CALCULS PAS À PAS :

x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
60	84	-8	-56	64	3136	448
62	95	-6	-45	36	2025	270
64	140	-4	0	16	0	0
66	155	-2	15	4	225	-30
68	119	0	-21	0	441	0
70	175	2	35	4	1225	70
72	145	4	5	16	25	20
74	197	6	57	36	3249	342
76	150	8	10	64	100	80
SOMME 612	1260			$SC_{xx} = 240$	$SC_{yy} = 10\,426$	$SC_{xy} = 1200$
$\bar{x} = 68$	$\bar{y} = 140$					

LA DROITE DE RÉGRESSION SE DÉTERMINE À PARTIR DES RÉSUMÉS STATISTIQUES EN BAS DU TABLEAU.

$$b = \frac{SC_{xy}}{SC_{xx}} = \frac{1200}{240} = 5 \quad a = \bar{y} - b\bar{x} = 140 - 5 \times 68 = -200$$

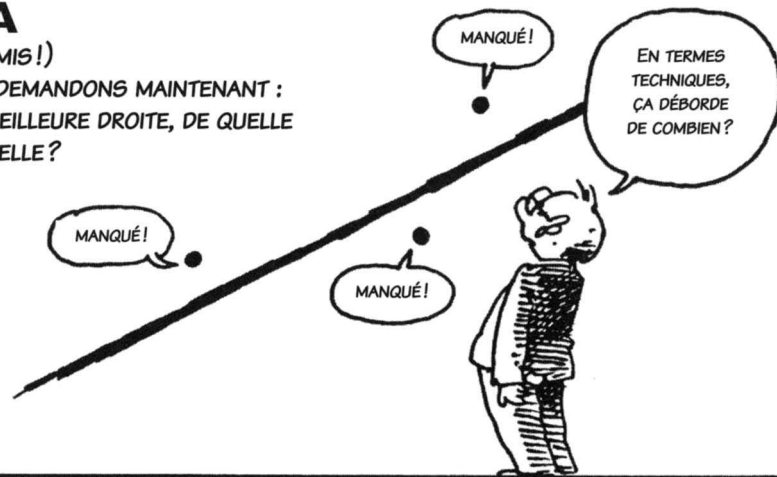
DONC $y = -200 + 5$



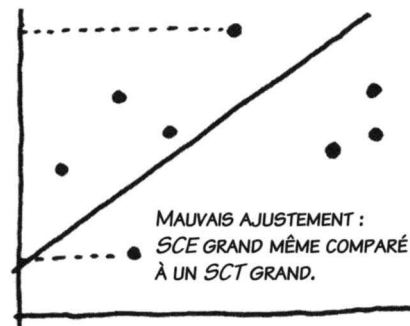
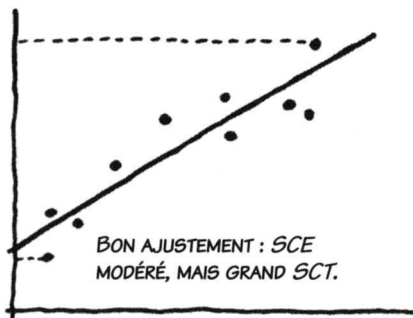
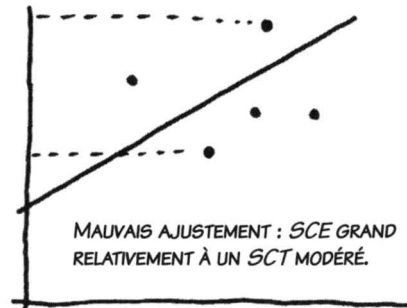
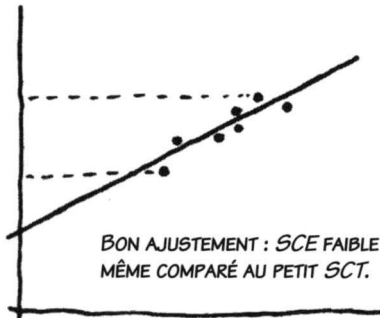
ANOVA

(COMME PROMIS!)

NOUS NOUS DEMANDONS MAINTENANT :
SI C'EST LA MEILLEURE DROITE, DE QUELLE
QUALITÉ EST-ELLE?



COMME VOUS POUVEZ L'IMAGINER, LA RÉPONSE À CETTE QUESTION DÉPEND DE LA FAÇON
DONT LES DONNÉES SONT ÉPARPILLÉES : DE L'IMPORTANCE DE **SCE**, LA SOMME
DES CARRÉS, RELATIVEMENT AUX **ÉCARTS TOTAUX** DES DONNÉES. QUELQUES EXEMPLES :



QUANTIFIONS CELA EN **SCINDANT**
LA VARIATION EN y . REGARDEZ
 LA FIGURE DE DROITE POUR SUIVRE.
 NOUS AVONS

$$\hat{y}_i = a + bx_i$$

OÙ \hat{y}_i EST LE POIDS **ESTIMÉ**
 PAR LA DROITE DE RÉGRESSION.

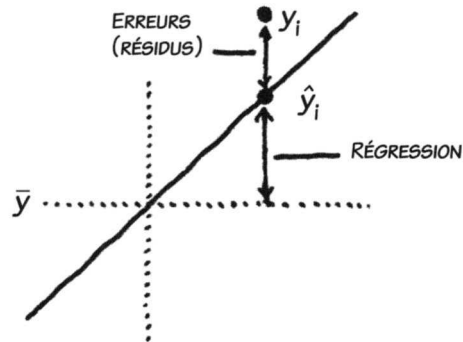


Table ANOVA

SOURCE DE VARIATION	SOMMES DES CARRÉS	VALEURS POUR LES DONNÉES
RÉGRESSION	$SCR = \sum (\hat{y}_i - \bar{y})^2$	6 000
ERREURS (RÉSIDUS)	$SCE = \sum (y_i - \hat{y}_i)^2$	4 426
TOTAL	$SCT = \sum (y_i - \bar{y})^2$	10 426

(À PROPOS, CE N'EST PAS ÉVIDENT QUE $SCT = SCR + SCE$, MAIS C'EST VRAI!)
 SINON, VOICI LE CALCUL DÉTAILLÉ DES SOMMES DE CARRÉS POUR LA RÉGRESSION
 ET LES RÉSIDUS POUR NOS DONNÉES AVEC $y = -200 + 5x$.

x_i	y_i	\hat{y}_i	RÉGRESSION		ERREURS (RÉSIDUS)	
			$(\hat{y}_i - \bar{y})$	$(\hat{y}_i - \bar{y})^2$	$(y_i - \hat{y}_i)$	$(y_i - \hat{y}_i)^2$
60	84	100	-40	1600	-16	256
62	95	110	-30	900	-15	225
64	140	120	-20	400	20	400
66	155	130	-10	100	25	625
68	119	140	0	0	-21	441
70	175	150	10	100	25	625
72	145	160	20	400	-15	225
74	197	170	30	900	27	729
76	150	180	40	1600	-30	900
$\bar{x} = 68$	$\bar{y} = 140$		$SCR = 6000$		$SCE = 4426$	

SCR MESURE LA VARIATION
TOTALE DUE À LA RÉGRESSION
(LES VALEURS PRÉDITES DE y).
SCE NOUS L'AVONS DÉJÀ VU.
NOTEZ QUE

$$\frac{SCE}{SCT}$$

EST LA PROPORTION
DES ERREURS (OU RÉSIDUS)
RELATIVEMENT AUX ÉCARTS
TOTAUX.

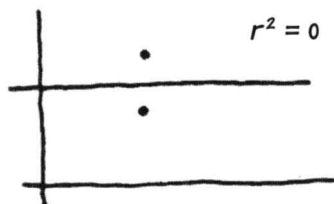


LE COEFFICIENT DE DÉTERMINATION

EST LA PROPORTION DE LA SCT IMPUTABLE
À LA RÉGRESSION :

$$r^2 = \frac{SCR}{SCT} = 1 - \frac{SCE}{SCT}$$

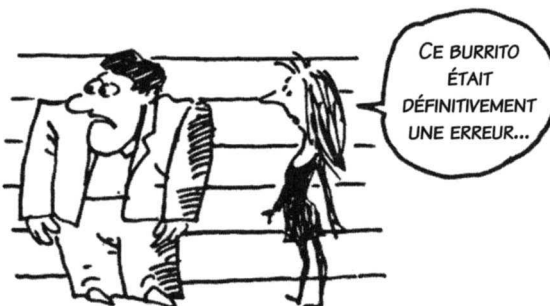
(CAR $SCR = SCT - SCE$). r^2 EST TOUJOURS
INFÉRIEUR OU ÉGAL À 1. PLUS C'EST PROCHE
DE 1, PLUS LES OBSERVATIONS SONT
PROCHES DE LA DROITE. $r^2 = 1$ INDIQUE
QUE LES POINTS SONT ALIGNÉS
SUR LA DROITE.



SI ON CALCULE, POUR NOTRE
PETIT JEU DE DONNÉES,
ON OBTIENT :

$$r^2 = \frac{6\,000}{10\,426} = 0,58$$

58 % DE LA VARIATION DU POIDS
S'EXPLIQUE PAR LA TAILLE.
LES 42 % RESTANTS SONT
UN TERME D'« ERREUR ».



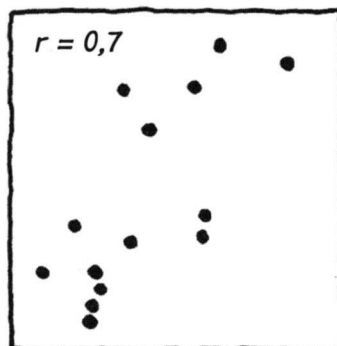
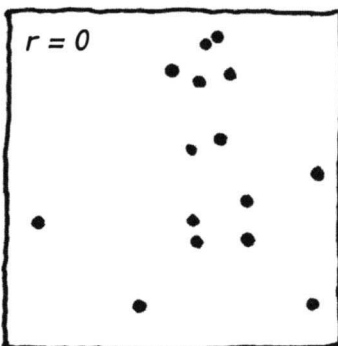
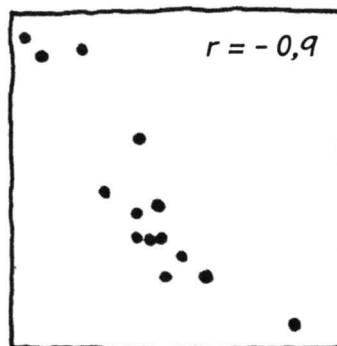
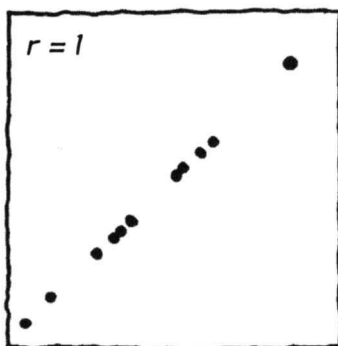
DE FAÇON ALTERNATIVE,
LE **coefficient de corrélation**
EST DÉFINI COMME LA RACINE
DE r^2 MULTIPLIÉE PAR LE SIGNE DE b .

$$r = (\text{signe de } b) \sqrt{r^2}$$

AINSI r EST POSITIF SI LA DROITE
EST CROISSANTE ET NÉGATIF
SI CELLE-CI EST DÉCROISSANTE.



r MESURE À LA FOIS LA JUSTESSE DE L'AJUSTEMENT ET LE SENS DE LA RELATION
ENTRE x ET y (EN INDICANT SI UN ACCROISSEMENT DE x AUGMENTE OU DIMINUE y).



MAINTENANT SOYONS
HONNÊTES : PERSONNE
OU PRESQUE NE
FAIT CES CALCULS
À LA MAIN DÉSORMAIS.
AVEC UN ORDINATEUR,
TOUT LE TRAVAIL PEUT
S'EFFECTUER AVEC
UNE LIGNE DE CODE.



EN FAIT, TOUT CE LIVRE
PEUT ÊTRE COMPRESSÉ
DANS LE CERVEAU
D'UN STATISTICIEN...

SI ON UTILISE LE LOGICIEL MINITAB, DÉVELOPPÉ À PENN STATE,
LA LIGNE DE COMMANDE EST :

MTB > régression 'poids' avec 1 variable indépendante 'taille'

ET LES RÉSULTATS SONT :

L'équation de la régression est

$$\text{Poids} = -200 + 5.00 * \text{taille}$$

Prédicteur	Coefficient	Écart-type	t Student	Valeur p
Constante	-200,0	110,7	-1,81	0,114
Taille	5,000	1,623	3,08	0,018

s = 25,15 R-carré = 57,5 % R-carré (ajusté) = 51,5 %

Analyse de Variance

SOURCE	DL	SC	MC	F	Valeur p
Régression	1	6 000,0	6 000,0	9,49	0,018
Erreur	7	4 426,0	632,3		
Total	8	10 426,0			

QUELLE CHARGE!



QUEL BONHEUR!
L'ORDINATEUR
EST D'ACCORD
AVEC NOUS!

MAINTENANT, FAISONS LA MÊME CHOSE AVEC LES VRAIES DONNÉES
DES 92 ÉTUDIANTS :

MTB > régression 'poids' avec 1 variable indépendante 'taille'

ET LES RÉSULTATS SONT :

L'équation de la régression est

$$\text{Poids} = -205 + 5,09 * \text{taille}$$

Prédicteur	Coefficient	Écart-type	t Student	Valeur p
Constante	-204,74	29,16	-7,02	0,000
Taille	5,0918	0,4237	12,02	0,000

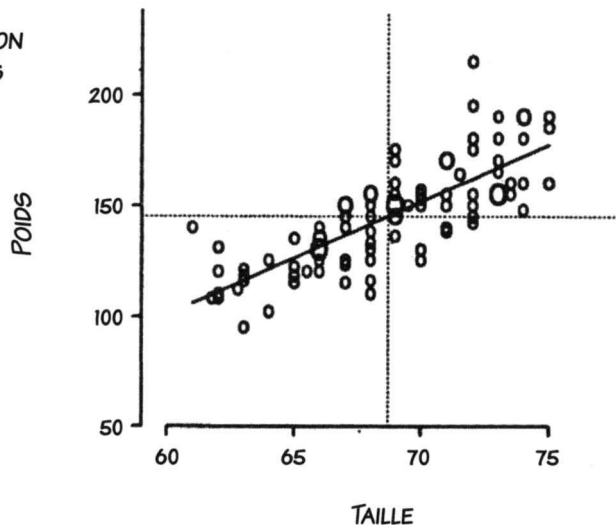
s = 14,79 R-carré = 61,6 % R-carré (ajusté) = 61,2 %

Analyse de Variance

SOURCE	DL	SC	MC	F	Valeur p
Régression	1	31 592	31 592	144,38	0,000
Erreur	7	19 692	219		
Total	8	51 284			

VOICI LE NUAGE DE POINTS
AVEC LA DROITE ESTIMÉE.
LE COEFFICIENT DE CORRÉLATION
POUR NOTRE JEU DE DONNÉES
VAUT

$$r = +\sqrt{0,616} = 0,78$$



INFÉRENCE STATISTIQUE

JUSQU'À MAINTENANT, NOUS AVONS FAIT DE L'ANALYSE DE DONNÉES, EN DÉTERMINANT LA MEILLEURE RELATION LINÉAIRE ENTRE LES DONNÉES OBSERVÉES DE x ET y . CHANGEONS MAINTENANT DE POINT DE VUE ET CONSIDÉRONS LES 92 ÉTUDIANTS COMME UN ÉCHANTILLON DE LA POPULATION GLOBALE DES ÉTUDIANTS. QUE POUVONS-NOUS EN DÉDUIRE ?



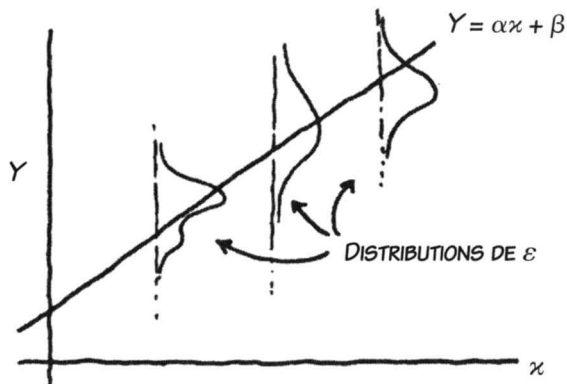
UN MODÈLE DE RÉGRESSION POUR LA POPULATION GLOBALE EST UNE RELATION LINÉAIRE

$$Y = \alpha + \beta x + \varepsilon$$

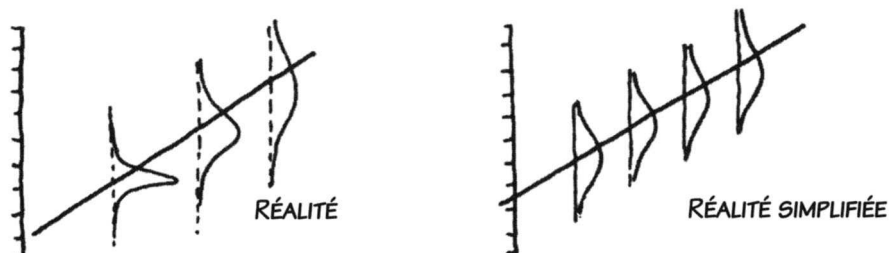
NOTEZ QUE
LES LETTRES GRECQUES
REPRÉSENTENT LES
PARAMÈTRES DU MODÈLE!

Y EST LA VARIABLE ALÉATOIRE DÉPENDANTE; x EST LA VARIABLE INDÉPENDANTE (QUI PEUT ÊTRE OU NON ALÉATOIRE); α ET β SONT LES PARAMÈTRES INCONNUS À ESTIMER; ET ε REPRÉSENTE LES ERREURS ALÉATOIRES FLUCTUANTES.

POUR LE MODÈLE
DE POIDS ET DE TAILLE,
 Y EST LE POIDS,
 x EST LA TAILLE,
 α ET β SONT INCONNUES.
ON PEUT VOIR ε COMME
UNE **COMPOSANTE**
ALÉATOIRE DES POIDS Y
POUR CHAQUE VALEUR
DE TAILLE x .



LA DISTRIBUTION DE ε N'EST EN FAIT PAS LA MÊME POUR DES VALEURS DIFFÉRENTES DE x . LES PERSONNES DE 1,50 m VARIENT MOINS EN POIDS QUE CELLES DE 1,80 m. NÉANMOINS, NOUS FAISONS MAINTENANT UNE HYPOTHÈSE SIMPLIFICATRICE : NOUS SUPPOSONS QUE POUR TOUTE VALEUR DE x LES ε SONT **INDÉPENDANTS ET NORMAUX**, QU'ILS ONT LE MÊME ÉCART-TYPE $\sigma = \sigma(\varepsilon)$ ET UNE MOYENNE $\mu = 0$.



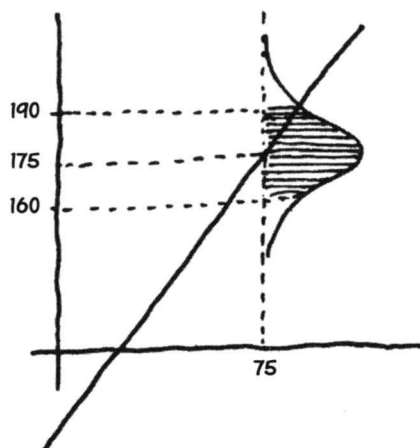
PAR EXEMPLE, SI LE VRAI MODÈLE DE POIDS EST LE SUIVANT :

$$Y = -125 + 4x + \varepsilon$$

OÙ ε SUIV UNE LOI NORMALE AVEC $\mu = 0$ ET $\sigma = 15$ LIVRES, ALORS, SELON CE MODÈLE, LES ÉTUDIANTS MESURANT 75 POUCES (1,90 m) ONT UNE DISTRIBUTION DE :

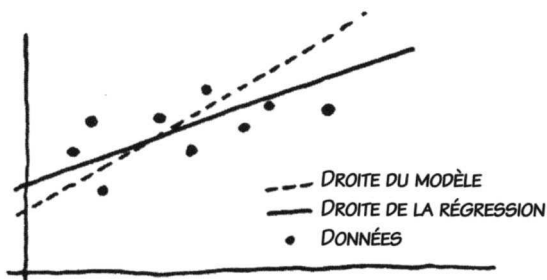
$$\begin{aligned} Y &= -125 + 4 \times 75 + \varepsilon \\ &= 175 + \varepsilon \end{aligned}$$

AINSI POUR $x = 75$, Y SUIV UNE LOI NORMALE DE MOYENNE 175 LIVRES ET D'ÉCART-TYPE 15 LIVRES.



MAINTENANT, ÉTANT DONNÉ NOTRE MODÈLE $Y = \alpha + \beta x + \varepsilon$, NOUS VOULONS, COMME NOUS L'AVONS FAIT DANS LES PRÉCÉDENTS CHAPITRES, PRENDRE UN **ÉCHANTILLON** ET L'UTILISER POUR **ESTIMER** α ET β .

ON PEUT MONTRER QUE α ET β , LES COEFFICIENTS DE LA DROITE DES MOINDRES CARRÉS $y = a + bx$ SONT **BLUE : BEST LINEAR UNBIASED ESTIMATORS**, DONC LES MEILLEURS ESTIMATEURS LINÉAIRES NON BIAISÉS DE α ET β .



GARANTIE
INCONDITIONNELLE!



COMME D'HABITUDE, DES ÉCHANTILLONS DIFFÉRENTS DONNANT LIEU À DES OBSERVATIONS DIFFÉRENTES NE GÉNÉRERAIENT PAS LA MÊME DROITE DE RÉGRESSION. CES DROITES SONT **DISTRIBUÉES** AUTOUR DE LA DROITE $Y = \alpha + \beta x + \varepsilon$. NOTRE QUESTION EST : COMMENT a ET b SONT-ILS DISTRIBUÉS AUTOUR DE α ET β RESPECTIVEMENT, ET COMMENT CONSTRUIRE DES **INTERVALLES DE CONFIANCE** ET DES **TESTS D'HYPOTHÈSES** ?

ILS SONT
BLUE...
ET JE SUIS VERT...

HEUREUSEMENT,
ÇA NE ME DÉRANGE
PAS D'ÊTRE VERT...

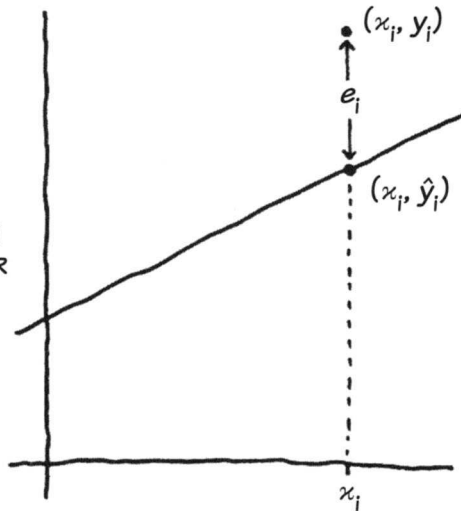


POUR CHAQUE OBSERVATION (x_i, y_i) ,
NOUS AVONS :

$$y_i = a + bx_i + e_i$$

OÙ $e_i = y_i - \hat{y}_i$ EST L'ÉCART ENTRE y_i
ET LA DROITE DE RÉGRESSION. LES e_i
SONT DES **VALEURS D'ÉCHANTILLON**
DE ε ET ELLES DONNENT UN ESTIMATEUR
 s POUR $\sigma(\varepsilon)$.

$$s = \sqrt{\frac{\sum e_i^2}{n-2}}$$



(ON MET $n - 2$ AU DÉNOMINATEUR CAR ON UTILISE DEUX DEGRÉS DE LIBERTÉ
POUR LE CALCUL DE a ET b , LAISSANT $n - 2$ INFORMATIONS INDÉPENDANTES
POUR ESTIMER σ .)

BIEN QUE CELA PUISSE PARAÎTRE
MOINS ÉVIDENT, ON PEUT AUSSI
ÉCRIRE s COMME CECI :

$$s = \sqrt{\frac{sc_{yy} - bsc_{xy}}{n-2}}$$

CETTE FORMULE NOUS PERMET
DE CALCULER s DIRECTEMENT
À PARTIR DES STATISTIQUES
D'ÉCHANTILLON.



RÉPÉTONS : s EST UN ESTIMATEUR DU DEGRÉ
D'ÉPARILLEMENT DES DONNÉES AUTOUR
DE LA DROITE.



INTERVALLES DE CONFIANCE

LES INTERVALLES DE CONFIANCE À 95 %
POUR α ET β ONT LA FORME FAMILIÈRE
(VOIR P. 133) SUIVANTE :

$$\beta \in b \pm t_{0,025} s(b)$$

$$\alpha \in a \pm t_{0,025} s(a)$$

OÙ ON UTILISE UN t DE STUDENT AVEC $n - 2$
DEGRÉS DE LIBERTÉ POUR LA MÊME RAISON
QUE PRÉCÉDEMMENT.



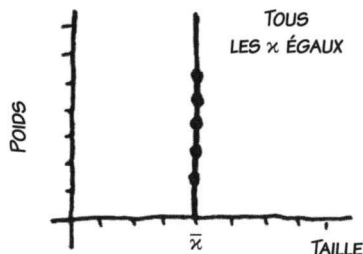
LES ERREURS-TYPES, PAR CONTRE, NE SONT PAS COMME D'HABITUDE.
VOICI LES FORMULES (SANS DÉMONSTRATIONS) :

$$s(b) = \frac{s}{\sqrt{sc_{xx}}}$$

$$s(a) = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{sc_{xx}}}$$



QU'EST-IL ARRIVÉ À NOTRE PRÉCIEUX $1/\sqrt{n}$? IL A ÉTÉ REMPLACÉ PAR sc_{xx} .
COMME n , sc_{xx} AUGMENTE AVEC LE NOMBRE DE DONNÉES, MAIS IL REFLÈTE AUSSI
LA VARIATION TOTALE DES DONNÉES EN x . PAR EXEMPLE, SI TOUS LES ÉTUDIANTS
DE L'ÉCHANTILLON AVAIENT LA MÊME TAILLE, NOUS NE POURRIONS EN TIRER AUCUNE
CONCLUSION SUR LA DÉPENDANCE DU POIDS PAR RAPPORT À LA TAILLE. DANS CE CAS,
 $sc_{xx} = 0$, CE QUI IMPLIQUERAIT $s(b) = \infty$ ET DES INTERVALLES DE CONFIANCE
DE LARGEUR INFINIE.



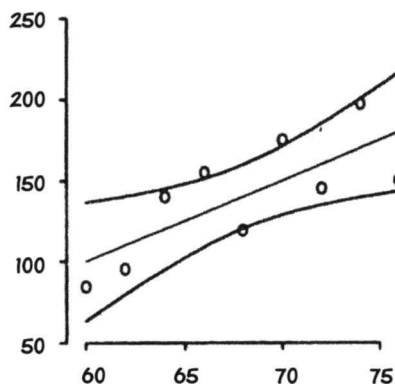
PLUS DE QUESTIONS :

COMMENT PRÉDIRE LA MOYENNE DES RÉPONSES Y POUR UNE VALEUR FIXÉE DE x_0 ? PAR EXEMPLE, QUEL EST LE POIDS MOYEN DES ÉTUDIANTS DE 76 POUCES? L'INTERVALLE DE CONFIANCE À 95 % EST :

$$\alpha + \beta x_0 \in a + bx_0 \pm t_{0,025} s(\hat{y})$$

où :

$$s(\hat{y}) = s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SC_{xx}}}$$



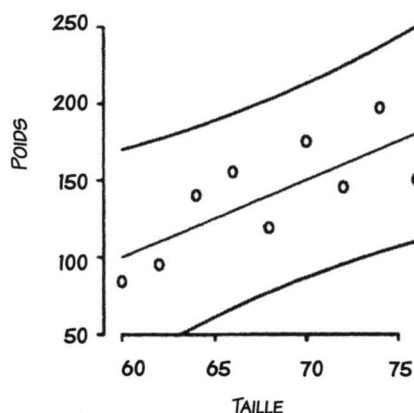
SUPPOSONS QU'UN NOUVEL ÉTUDIANT DE TAILLE x_{nv} ARRIVE. QUELLE PRÉVISION POUVONS-NOUS FAIRE SUR SON POIDS y_{nv} ?

L'INTERVALLE DE PRÉDICTION À 95 % DE y_{nv} POUR L'INDIVIDU DE TAILLE x_{nv} EST :

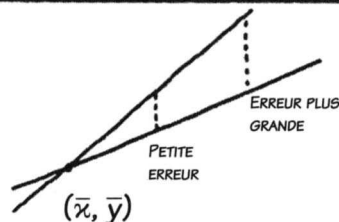
$$y_{nv} \in a + bx_{nv} \pm t_{0,025} s(y_{nv})$$

où

$$s(y_{nv}) = s \sqrt{1 + \frac{1}{n} + \frac{(x_{nv} - \bar{x})^2}{SC_{xx}}}$$



CES DEUX ERREURS-TYPES CONTIENNENT UN TERME CROISSANT EN x (RESPECTIVEMENT EN x_0 ET x_{nv}) QUAND ON S'ÉLOIGNE DE LA MOYENNE \bar{x} . POURQUOI L'ERREUR AUGMENTE-T-ELLE EN S'ÉLOIGNANT DE \bar{x} ? PARCE QUE COMME LA DROITE PASSE TOUJOURS PAR LE POINT (\bar{x}, \bar{y}) , SI L'ON FAIT PIVOTER CETTE DROITE DE RÉGRESSION ALORS L'IMPACT EST PLUS ÉLEVÉ SUR LES VALEURS ÉLOIGNÉES DE LA MOYENNE!



REGARDONS CELA AVEC NOS DONNÉES ARRANGÉES.
 POUR LA MOYENNE DES POIDS POUR UNE TAILLE
 $x = 76$ POUCES, NOUS AVONS $b = -200$ ET $a = 5$.
 DONC :

$$Y \in -200 + 5 \times 76 \pm 2,365 \times 25,15 \sqrt{\frac{1}{q} + \frac{(76 - 68)^2}{240}}$$

$$Y \in 180 \pm 2,365 \times 25,15 \sqrt{0,3777}$$

$$Y \in 180 \pm 36,55 \text{ LIVRES}$$

$$Y \in 81,6 \pm 16,6 \text{ kg}$$

LA MOYENNE ESTIMÉE DU POIDS DES ÉTUDIANTS
 DE 1,93 m EST DE 81,6 kg, ET NOUS SOMMES
 SÛRS À 95 % QU'ELLE NE DÉVIERA PAS DE PLUS
 DE 16,6 kg.



SI L'ON PREND UN NOUVEL ÉTUDIANT DE 76 POUCES, NOUS UTILISONS NOS DONNÉES
 ARRANGÉES DE NEUF OBSERVATIONS POUR PRÉDIRE QUE

$$y_{nv} \in -200 + 5 \times 76 \pm 2,365 \times 25,15 \sqrt{1 + \frac{1}{q} + \frac{(76 - 68)^2}{240}}$$

$$y_{nv} \in 180 \pm 2,365 \times 25,15 \times 1,174$$

$$y_{nv} \in 180 \pm 70 \text{ LIVRES}$$

$$y_{nv} \in 81,6 \pm 31,7 \text{ kg}$$

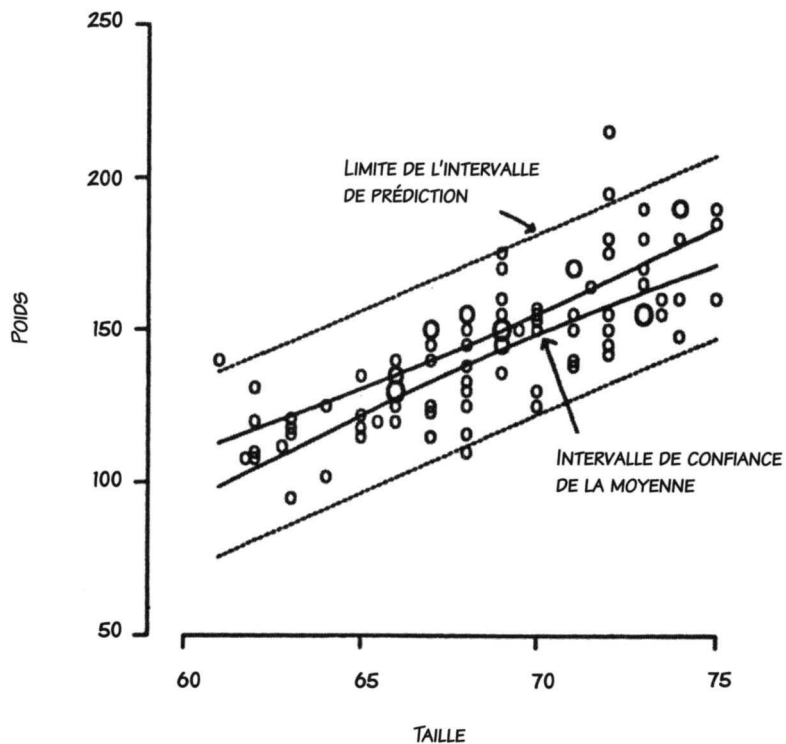


ON PEUT DIRE
 À L'ENTRAÎNEUR
 DE FOOT QUE
 NOUS SOMMES
 PRESQUE SÛRS
 QUE CE NOUVEL
 ÉTUDIANT PÈSE
 ENTRE 50
 ET 113 kg!

LES INTERVALLES NE SONT PAS FAMEUX! QUEL EST LE PROBLÈME?
EN FAIT, IL Y A DEUX PROBLÈMES.

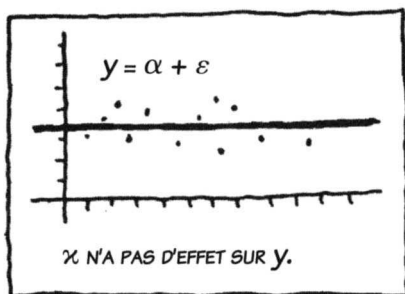


L'ENSEMBLE DES ÉTUDIANTS DE PENN STATE DONNE DE MEILLEURS RÉSULTATS.



Tests d'hypothèses

QUELQU'UN DE COMPLÈTEMENT SCEPTIQUE POURRAIT SUGGÉRER QU'IL N'Y A PAS DE RELATION ENTRE LA TAILLE ET LE POIDS. CE QUI REVIENT À DIRE QUE $\beta = 0$.



NOUS PRENONS CECI COMME NOTRE **HYPOTHÈSE NULLE**.

$$H_0 : \beta = 0$$

DANS CE CAS, LA STATISTIQUE DE TEST EST :

$$t = \frac{b}{s(b)}$$

C'EST UN t DE STUDENT AVEC $n - 2$ DEGRÉS DE LIBERTÉ. COMME D'HABITUDE, LA SIGNIFICATION DU TEST DÉPEND DE L'HYPOTHÈSE ALTERNATIVE :

$$t > t_{\alpha} \text{ POUR } H_a : \beta > 0$$

$$t < t_{\alpha} \text{ POUR } H_a : \beta < 0$$

$$|t| > |t_{\alpha/2}| \text{ POUR } H_a : \beta \neq 0$$

POUR LES DONNÉES ARRANGÉES DE POIDS, NOUS SUSPECTONS FORTEMENT QUE L'HYPOTHÈSE ALTERNATIVE DOIT ÊTRE :

$$H_a : \beta > 0$$

ON TESTE ALORS :

$$t_{obs} = \frac{5}{s(b)} = \frac{5}{1,62}$$

$$= 3,08$$

AVEC 7 DEGRÉS DE LIBERTÉ, $t_{0,05} = 1,895$. COMME $t_{obs} > t_{0,05}$, ON REJETTE L'HYPOTHÈSE NULLE AU SEUIL DE SIGNIFICATION DE 5 %, ET ON CONCLUT QU'IL Y A UNE RELATION POSITIVE SIGNIFICATIVE ENTRE LA TAILLE ET LE POIDS.



Régression linéaire multiple

LE MÊME TYPE D'IDÉES BASIQUES PEUT ÊTRE UTILISÉ POUR ANALYSER LES RELATIONS ENTRE UNE VARIABLE DÉPENDANTE ET **PLUSIEURS** VARIABLES INDÉPENDANTES.

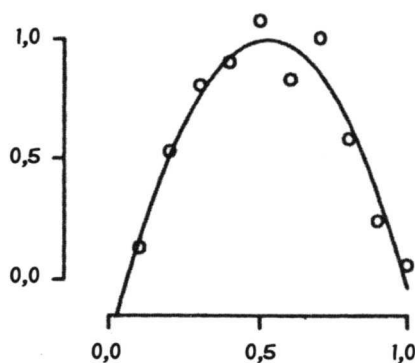
$$Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

PAR EXEMPLE, LE POIDS EST DÉTERMINÉ PAR UN NOMBRE DE FACTEURS AUTRES QUE LA TAILLE COMME L'ÂGE, LE SEXE, LES RÉGIMES, LE TYPE PHYSIQUE, ETC.



L'ALGÈBRE MATRICIELLE ET LES ORDINATEURS SE COMBINENT POUR FACILITER L'ANALYSE DE CES PROBLÈMES.

Régression non linéaire



PARFOIS LES DONNÉES SUIVENT DE FAÇON ÉVIDENTE UNE COURBE **NON LINÉAIRE**. LES STATISTICIENS ONT DES TAS D'ASTUCES POUR UTILISER DES TECHNIQUES DE RÉGRESSIONS **LINÉAIRES** POUR DES PROBLÈMES **NON LINÉAIRES**. LA PLUS SIMPLE EST DE CONSIDÉRER Y COMME UN POLYNÔME :

$$Y = \alpha + \beta_1 x + \beta_2 x^2 + \varepsilon$$

ET DE TRAITER x ET x^2 COMME DES VARIABLES INDÉPENDANTES D'UN MODÈLE LINÉAIRE.

Diagnostic des régressions

AJUSTER UN MODÈLE COMPLEXE À DES DONNÉES PEUT PARFOIS OCCULTER DES DIFFICULTÉS OU DES PROBLÈMES. ON UTILISE DES PROCÉDURES DE DIAGNOSTIC DE RÉGRESSION POUR RÉVÉLER TOUTES LES SURPRISES CACHÉES DÉSAGRÉABLES.

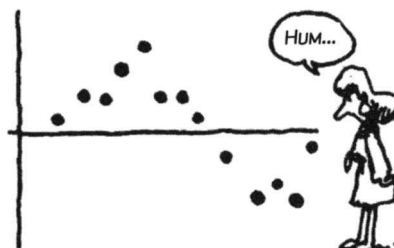


LA MÉTHODE LA PLUS SIMPLE EST DE DESSINER LES ERREURS OU **RÉSIDUS** e_i EN FONCTION DE LA VARIABLE **EXPLICATIVE** y_i . RAPPELEZ-VOUS QUE L'ERREUR ε ÉTAIT CENSÉE ÊTRE INDÉPENDANTE DE x .

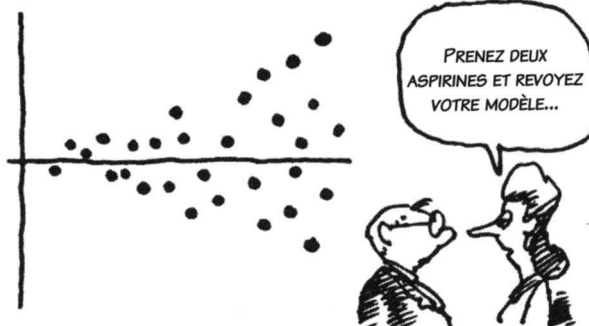
UN NUAGE **ALÉATOIRE** DE POINTS INDIQUE QUE LES HYPOTHÈSES DU MODÈLE SONT SÛREMENT CORRECTES.



TOUT AUTRE **SCHEMA** INDIQUE UN VÉRITABLE PROBLÈME AVEC LES HYPOTHÈSES DU MODÈLE.



UNE SURPRISE DÉSAGRÉABLE CLASSIQUE (QUE L'ON RETROUVE DANS NOS DONNÉES DE TAILLE/POIDS) EST QUE LES ERREURS SOIENT **HÉTÉROSCÉDASTIQUES**, QUAND LA VARIATION DE e CROÎT LORSQUE y AUGMENTE.



DANS CE CHAPITRE, NOUS AVONS
RÉSUMÉ LES IDÉES DE BASE
ET LES TECHNIQUES DE L'ANALYSE
DE RÉGRESSION, QUI ÉTUDIE
LES RELATIONS STATISTIQUES
ENTRE DES VARIABLES. AINSI S'ACHÈVE
NOTRE DISCUSSION DÉTAILLÉE SUR
LES MÉTHODES STATISTIQUES DE BASE.
DANS LE DERNIER CHAPITRE, NOUS
EXAMINERONS BRIÈVEMENT QUELQUES
AUTRES SUJETS ET PROBLÈMES.



Chapitre 12

CONCLUSION

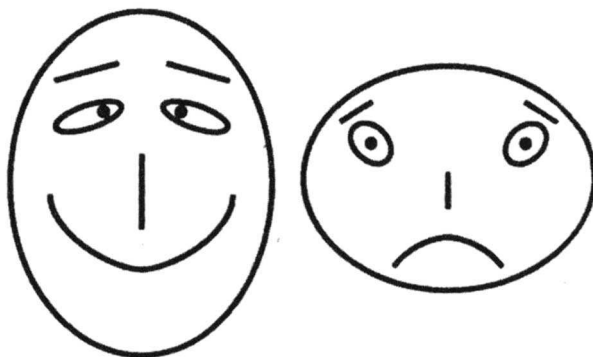
LES PRINCIPES DE BASE, LES OUTILS ET LES CALCULS
ÉTUDIÉS DANS CE LIVRE PEUVENT ÊTRE PLUS LARGEMENT
UTILISÉS POUR RÉSOUDRE DES PROBLÈMES PLUS COMPLEXES.
VOICI UN **ÉCHANTILLON BIAISÉ** DE MÉTHODES
STATISTIQUES PLUS AVANCÉES !



AFFICHAGE DE DONNÉES

α = INCLINAISON DES SOURCILS
 y = TAILLE DES YEUX
 z = TAILLE DU NEZ

t = TAILLE DE LA BOUCHE
 β = HAUTEUR DU VISAGE
 ETC.



ON A VU COMMENT REPRÉSENTER UNE VARIABLE AVEC UN GRAPHE DE POINTS ET DEUX VARIABLES AVEC UN NUAGE DE POINTS. MAIS COMMENT REPRÉSENTER GRAPHIQUEMENT **PLUS DE DEUX VARIABLES** SUR UNE PAGE PLANE ?
 PARMIS LES POSSIBILITÉS, UN GUIDE EN BD SE DOIT DE MENTIONNER L'IDÉE SIMPLE DE **HERMAN CHERNOFF** (1923-) : UTILISER UN VISAGE HUMAIN OÙ CHAQUE TRAIT CORRESPOND À UNE VARIABLE POUR CRÉER UN **VISAGE DE CHERNOFF**.

Analyse statistique de DONNÉES MULTIVARIÉES

IL EXISTE DES MODÈLES MULTIVARIÉS POUR AIDER À L'ANALYSE ET À L'AFFICHAGE DE DONNÉES À n DIMENSIONS. VOICI QUELQUES TECHNIQUES MULTIVARIÉES :

L'analyse typologique

VOISE À DIVISER LA POPULATION EN SOUS-GROUPES HOMOGÈNES. PAR EXEMPLE, QUAND ON ANALYSE LES SCHÉMAS DE VOTE DU CONGRÈS, ON PEUT VOIR QUE LES REPRÉSENTANTS DU SUD-EST ET CEUX DU NORD-EST FORMENT DEUX GROUPES DISTINCTS.



L'analyse discriminante

EST LE PROCÉDÉ INVERSE. PAR EXEMPLE, LE BUREAU D'ADMISSION EN LICENCE AIMERAIT DISPOSER DE DONNÉES QUI LUI INDIQUERAIENT PAR AVANCE SI LE CANDIDAT **RÉUSSIRA** SON DIPLÔME (CONTRIBUANT ALORS AU FONDS DES ANCIENS) OU S'IL SERA EN **ÉCHEC** (QUITTANT L'UNIVERSITÉ POUR UNE RECONVERSION DANS L'HUMANITAIRE).



L'analyse factorielle

VISE À EXPLIQUER DES DONNÉES À DIMENSIONS MULTIPLES AVEC UN PETIT NOMBRE DE VARIABLES. UN PSYCHOLOGUE PEUT PAR EXEMPLE PROPOSER UN TEST DE 100 QUESTIONS EN PENSANT QUE LES RÉPONSES NE DÉPENDENT QUE DE QUELQUES FACTEURS COMME L'EXTRAVERSION, L'AUTORITARISME ET L'ALTRUISME. ON PEUT ALORS RÉSUMER LES RÉSULTATS DU TEST À PARTIR DE SCORES COMBINÉS DANS CES DIMENSIONS.

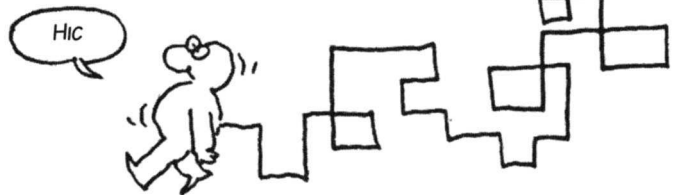


IL Y A AUSSI PLUS À FAIRE AVEC LES

PROBABILITÉS :

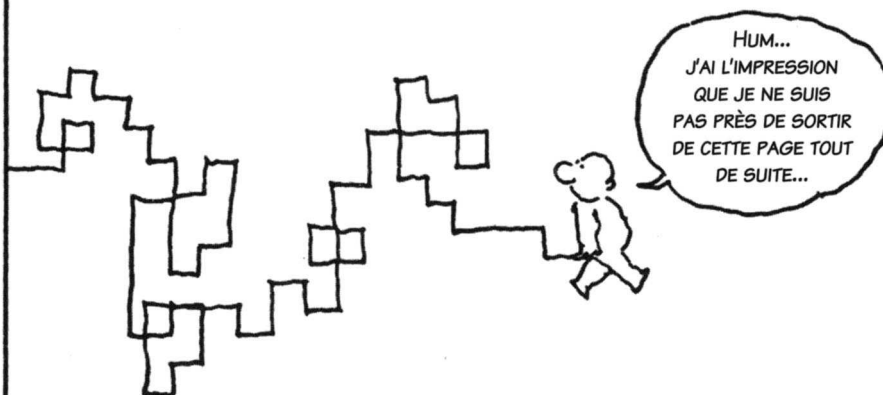
Une marche aléatoire

ELLE COMMENCE PAR UN LANCER DE PIÈCE.
SUPPOSONS QUE VOUS AVANCIEZ SI C'EST FACE
ET QUE VOUS RECLIEZ SI C'EST PILE (AVEC DEUX
PIÈCES ON PEUT FAIRE CELA SUR DEUX DIMENSIONS).
CES LANCERS RÉPÉTÉS PRODUISENT UN PROCESSUS
STOCHASTIQUE APPELÉ MARCHÉ ALÉATOIRE.
LES MODÈLES DE MARCHÉ ALÉATOIRE SONT
UTILISÉS POUR LES NÉGOCIATIONS EN BOURSE
ET LA GESTION DE PORTEFEUILLE.



L'analyse de séries temporelles

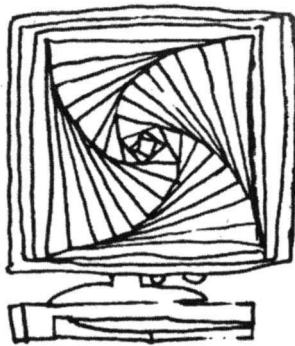
ELLE TRAITE D'ENSEMBLES DE DONNÉES, COMME LES MARCHES ALÉATOIRES, QUI VARIENT
AU COURS DU TEMPS : LES TEMPÉRATURES GLOBALES ET LOCALES, LE PRIX DU PÉTROLE,
ETC. L'ANALYSE DE SÉRIES TEMPORELLES SE FAIT À L'AIDE DE MODÈLES ALÉATOIRES
AFIN DE PRÉDIRE LES VALEURS FUTURES.



NOUS AVONS DÉJÀ VU L'AIDE APPORTÉE PAR LES ORDINATEURS POUR L'ARITHMÉTIQUE ET L'ANALYSE. IL Y A AUSSI DES IDÉES STATISTIQUES DONT L'**EXISTENCE** MÊME PROVIENT DES ORDINATEURS.

Analyse d'image

UNE IMAGE D'ORDINATEUR EST CONSTITUÉE DE PLUSIEURS MILLIONS DE PIXELS (ÉLÉMENTS DE L'IMAGE). UNE IMAGE DONNÉE EST DONC UN ENSEMBLE HAUTEMENT STRUCTURÉ DE MILLIONS DE VECTEURS DE PIXELS. L'ANALYSE D'IMAGES TEND À EXTRAIRE DU SENS DE CE TYPE D'INFORMATION.



ON UTILISE DES IMAGES POUR
COMPRENDRE LES DONNÉES,
MAIS MAINTENANT NOUS DEVONS
COMPRENDRE LES IMAGES!

Ré-échantillonnage

PARFOIS, LES ÉCARTS-TYPES ET LES BORNES DE CONFIANCE SONT IMPOSSIBLES À DÉTERMINER. ON UTILISE ALORS LE RÉ-ÉCHANTILLONNAGE, UNE TECHNIQUE QUI SE SERT DE L'ÉCHANTILLON LUI-MÊME **COMME S'IL S'AGISSAIT D'UNE POPULATION**. LES TECHNIQUES S'APPELLENT **RANDOMISATION, JACKKNIFE, ET TECHNIQUE DU BOOTSTRAPPING**.



NGH! ÇA A L'AIR
IMPOSSIBLE,
MAIS ÇA FONCTIONNE!

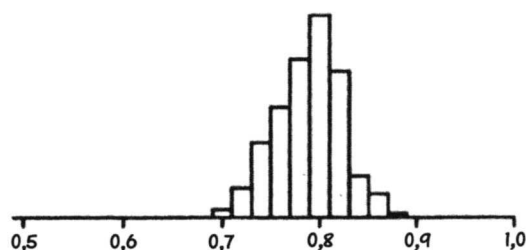
Ré-échantillonnage (suite)

POUR RÉ-ÉCHANTILLONNER, L'ORDINATEUR

- RE-ÉCHANTILLONNE L'ÉCHANTILLON;
- CALCULE LES ESTIMATIONS DE CES NOUVEAUX ÉCHANTILLONS;
- RÉPÈTE PLUSIEURS FOIS LES DEUX PREMIÈRES ÉTAPES, ET ESTIME LA VARIANCE DES ESTIMATIONS DES ÉCHANTILLONS DE L'ÉCHANTILLON.



VOUS VOUS SOUVENEZ DU COEFFICIENT DE CORRÉLATION r DES 92 POIDS ET TAILLES D'ÉTUDIANTS DU CHAPITRE 11 ? QUEL EST L'ÉCART-TYPE DE r ? L'ORDINATEUR PREND 200 ÉCHANTILLONS (« BOOTSTRAP ») DE CES 92 OBSERVATIONS, CALCULE r CHAQUE FOIS ET TRACE UN HISTOGRAMME DES r TROUVÉS.



CORRÉLATION DE BOOTSTRAPS



REMARQUEZ QUE LA VARIANCE DES ESTIMATIONS DES BOOTSTRAPS EST RELATIVEMENT FAIBLE.

ET FINALEMENT, VOICI
QUELQUES PROBLÈMES
À GARDER À L'ESPRIT...



LA QUALITÉ DES DONNÉES

VRAISEMBLABLEMENT DES PETITES ERREURS DANS L'ÉCHANTILLONNAGE, DANS LES MESURES ET L'ENREGISTREMENT DES DONNÉES PEUVENT CAUSER DES RAVAGES SUR N'IMPORTE QUELLE ANALYSE. R. A. FISHER (1890-1962), GÉNÉTICIEN ET FONDATEUR DES STATISTIQUES MODERNES, NE SE CONTENTAIT PAS DE CONCEVOIR ET D'ANALYSER DES EXPÉRIENCES CONCERNANT L'ÉLEVAGE DES ANIMAUX. IL NETTOYAIT AUSSI LEURS CAGES ET LES SOIGNAIT CAR IL SAVAIT QUE LA PERTE DE L'UN D'ENTRE EUX RISQUAIT D'AFPECTER SES RÉSULTATS.



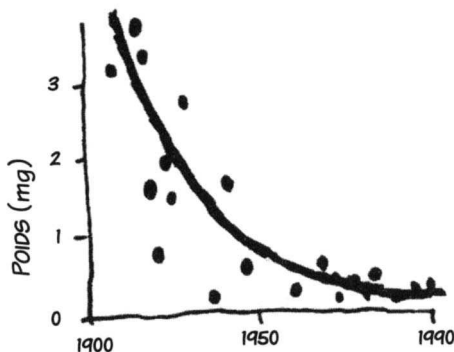
AVEC LEURS ORDINATEURS, LEURS BASES DE DONNÉES ET LEURS SUBVENTIONS GOUVERNEMENTALES, LES STATISTICIENS MODERNES ONT PERDU L'HABITUDE DE METTRE LA MAIN À LA PÂTE.



HÉ, JE SUIS GENTIL AVEC MA SOURIS MOI AUSSI!



SI VOUS TRACIEZ LA MASSE DE RÉSIDUS DE RATS SOUS LES ONGLES DES STATISTICIENS AU COURS DU SIÈCLE DERNIER, VOUS AURIEZ SÛREMENT UN GRAPHIQUE DE CE TYPE :



Innovation

LES MEILLEURES SOLUTIONS NE SONT PAS TOUJOURS DANS LES LIVRES !
PAR EXEMPLE, UNE COMPAGNIE RECRUTÉE POUR ESTIMER LA COMPOSITION
D'UNE **DÉCHARGE PUBLIQUE** A ÉTÉ CONFRONTÉE À DES PROBLÈMES INTÉRESSANTS
MAIS QUE VOUS NE TROUVEREZ PAS DANS LES TEXTES STANDARDS...



Communication

LES ANALYSES BRILLANTES SONT INUTILES SI LES RÉSULTATS NE SONT PAS CLAIREMENT EXPRIMÉS AVEC DES MOTS SIMPLES, Y COMPRIS SUR LE DEGRÉ D'INCERTITUDE DES CONCLUSIONS. PAR EXEMPLE, LES MÉDIAS PARLENT MAINTENANT DAVANTAGE DES MARGES D'ERREUR DE LEURS SONDAGES.



Travail d'équipe

DANS NOTRE SOCIÉTÉ COMPLEXE, LA RÉOLUTION DE NOMBREUX PROBLÈMES REQUIERT UN TRAVAIL D'ÉQUIPE. LES INGÉNIEURS, LES STATISTICIENS ET LES OUVRIERS À LA CHAÎNE COOPÈRENT POUR ACCROÎTRE LA QUALITÉ DE LEURS PRODUITS. LES BIOSTATISTIENS, LES MÉDECINS ET LES MILITANTS ONT TRAVAILLÉ ENSEMBLE POUR CONCEVOIR DES ESSAIS CLINIQUES AFIN D'ÉVALUER RAPIDEMENT L'EFFICACITÉ DE MÉDICAMENTS DESTINÉS À COMBATTRE LE SIDA.



EH BIEN VOILÀ ! MAINTENANT, VOUS DEVRIEZ POUVOIR FAIRE À PEU PRÈS N'IMPORTE QUOI
AVEC LES STATISTIQUES, EXCEPTÉ **MENTIR, TRICHER, VOLER** ET PARIER OU JOUER.

NOUS GARDONS
CES SUJETS POUR
LA BIBLIOGRAPHIE !



AVEZ-VOUS UNE ASSURANCE
CONVENABLE POUR FAUTE
PROFESSIONNELLE STATISTIQUE ?



BIBLIOGRAPHIE

POUR BIEN DÉBUTER

HAHN C., MACÉ S., *MÉTHODES STATISTIQUES APPLIQUÉES AU MANAGEMENT*, PARIS, PEARSON, 2012. PRÉSENTATION DES MÉTHODES STATISTIQUES SANS TROP DE FORMALISME ET PLUTÔT AXÉE SUR DES EXEMPLES PRATIQUES.

VESSEREAU A., *LA STATISTIQUE*, PARIS, PUF, COLL. « QUE SAIS-JE ? », 2002. UN PEU DE STATISTIQUES DANS UN FORMAT DE POCHÉ.

POUR APPROFONDIR ET POUR L'ÉTUDIANT

ANDERSON D. R., CAMM J. D., COCHRAN J. J. ET SWEENEY. D. J., *STATISTIQUES POUR L'ÉCONOMIE ET LA GESTION*, 5^e ÉDITION, DE BOECK, 2015. UN MANUEL TRÈS COMPLET AVEC DE NOMBREUX EXEMPLES ET EXERCICES.

LECOUTRE J.-P., *STATISTIQUE ET PROBABILITÉS*, 2^e ÉDITION, PARIS, DUNOD, 2003. COURS AVEC EXERCICES CORRIGÉS D'UN NIVEAU DE DEUXIÈME ANNÉE DE LICENCE D'ÉCONOMIE-GESTION.

PHAN T., ROWENCZYK J.-P., *EXERCICES ET PROBLÈMES DE STATISTIQUE ET PROBABILITÉS*, 2^e ÉDITION, PARIS, DUNOD, 2007. UN AUTRE COURS AVEC EXERCICES ET PROBLÈMES CORRIGÉS S'ADRESSANT À UN PUBLIC PLUS SCIENTIFIQUE QUE LE PRÉCÉDENT.

EN ANGLAIS

POUR L'ÉTUDIANT

FREEDMAN D., PISANI R. ET PURVES R., *STATISTICS*, 4th EDITION, NEW YORK, W. W. NORTON, 2007. INTRODUCTION APPROFONDIE AUX STATISTIQUES UTILISANT LE MOINS POSSIBLE LES MATHÉMATIQUES ET LES ÉQUATIONS.

MOORE D. S., NOTZ W. I., *STATISTICS CONCEPTS AND CONTROVERSIES*, NEW YORK, W. H. FREEMAN, 2008. SOULIGNE LES IDÉES PLUTÔT QUE LE CÔTÉ MÉCANIQUE.

CES DEUX LIVRES SONT DES CLASSIQUES : BONS, LITTÉRAIRES ET SPIRITUELS. À CÔTÉ DE CEUX-CI, IL EXISTE DES CENTAINES DE LIVRES DISPONIBLES ET LA PLUPART D'ENTRE EUX SONT, SELON NOUS, SATISFAISANTS.



COMMENT MENTIR, TRICHER ET PARIER :

VOS AUTEURS, QUI SONT DES SAINTS, MANQUENT D'EXPÉRIENCE EN LA MATIÈRE. VOICI QUELQUES CONSEILS DE PROS :

HUFF D., *HOW TO LIE WITH STATISTICS*, NEW YORK, W. W. NORTON, 1954. BON MARCHÉ ET TOUJOURS PUBLIÉ.

ORKIN M., *CAN YOU WIN?*, NEW YORK, W. H. FREEMAN, 1991. CONSEILS D'UN EXPERT DES PROBABILITÉS ET DES JEUX (N'EST PLUS PUBLIÉ, MAIS ON TROUVE ENCORE FACILEMENT DES EXEMPLAIRES D'OCCASION).

SPIRER H. F., SPIRER L. ET JAFFE A. J., *MISUSED STATISTICS*, 2nd EDITION, NEW YORK, MARCEL DECKER, 1998. UN LIVRE FAISANT PARTIE D'UNE COLLECTION RECONNUE SUR LES STATISTIQUES.

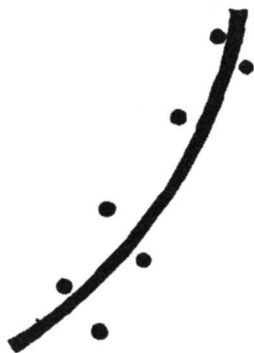


DROIT ET SOCIÉTÉ

GASTWIRTH J. L., *STATISTICAL REASONING IN LAW AND POLICY*, VOL. 1 ET 2, SAN DIEGO, ACADEMIC PRESS, 1988. LES DÉTAILS PRATIQUES DE LA LOI INCLUANT LE CAS DE SÉLECTION DE JURY DÉCRIT AU DÉBUT DU CHAPITRE 9.

DANS LE CHAPITRE 9, LE COMMENTAIRE NON JUDICIAIRE SUR LE POKER EST TIRÉ D'UNE AFFAIRE RÉELLE; NOUS EN AVONS EU LA CONFIRMATION PAR LE DR JOHN DE CANI DE L'UNIVERSITÉ DE PENNSYLVANIE.

LE COMITÉ DE PILOTAGE DU GROUPE DE RECHERCHE ET D'ÉTUDE SUR LA SANTÉ DES MÉDECINS, « FINAL REPORT ON THE ASPIRIN COMPONENT OF THE ONGOING PHYSICIANS' HEALTH STUDY », *THE NEW ENGLAND JOURNAL OF MEDICINE*, VOL. 321, P. 129-135. VOUS Y TROUVEREZ DES DÉTAILS CONCERNANT L'ÉTUDE SUR L'ASPIRINE DÉCRITE DANS LE CHAPITRE 9.



DESCRIPTION GRAPHIQUE DES DONNÉES

TUFTE E. R., *THE VISUAL DISPLAY OF QUANTITATIVE INFORMATION*, NEW HAVEN, GRAPHICS PRESS, 2001. LES LIVRES DE TUFTE ONT FIXÉ DE NOUVEAUX STANDARDS POUR LA COMMUNICATION DES DONNÉES QUANTITATIVES.

CLEVELAND W. S., *THE ELEMENTS OF GRAPHING DATA*, SUMMIT (N. J.), HOBART PRESS, 1994. LA PLUPART DES IDÉES GRAPHIQUES DÉCRITES ICI PEUVENT ÊTRE IMPLÉMENTÉES AVEC DES LOGICIELS DE STATISTIQUES COMME R OU D'AUTRES SYSTÈMES.

HISTOIRE

BOX J. F., R. A. FISHER : *THE LIFE OF A SCIENTIST*, NEW YORK, WILEY, 1978. LA BIOGRAPHIE FAITE PAR LA FILLE DE L'UN DES STATISTICIENS LES PLUS INFLUENTS ET CONTROVERSÉS DU XX^e SIÈCLE (ELLE N'EST PLUS PUBLIÉE ET EST DEVENUE PLUTÔT RARE ET CHÈRE).

DAVID F. N., *GAMES, GODS AND GAMBLING*, MINEOLA (N. Y.), DOVER 2012. RÉIMPRESSION D'UN LIVRE CLASSIQUE SUR L'HISTOIRE DES STATISTIQUES.

SALSBURG D., *THE LADY TASTING TEA : HOW STATISTICS REVOLUTIONIZED SCIENCE IN THE TWENTIETH CENTURY*, NEW YORK, HENRY HOLT & CO., 2002. UNE HISTOIRE POPULAIRE DES STATISTIQUES MODERNES.



LOGICIEL DE STATISTIQUES

DANS CE LIVRE, NOUS AVONS UTILISÉ LE LOGICIEL STATISTIQUE MINITAB (INC. STATE COLLEGE PA). LES DONNÉES DE POIDS ET DE TAILLES DES ÉTUDIANTS DE PENN STATE PROVIENNENT DE LA BASE DE DONNÉES PULSE DE CE SYSTÈME. LA VERSION ACTUELLE POUR ÉTUDIANTS DE MINITAB EST MINITAB EXPRESS. LES FICHIERS D'AIDE DE MINITAB SONT CLAIRS ET DÉTAILLÉS, IL S'AGIT PRESQUE EN SOI D'UN COURS BASIQUE DE STATISTIQUES. LES GRAPHIQUES INFORMATIQUES DE CETTE BD ONT ÉTÉ GÉNÉRÉS AVEC S-PLUS.

R EST UN SYSTÈME SIMILAIRE ET GRATUIT DE STATISTIQUES, TRÈS LARGEMENT UTILISÉ POUR DES ANALYSES GRAPHIQUES OU NUMÉRIQUES DÉTAILLÉES. R ET S ONT, TOUS LES DEUX, ÉTÉ DÉVELOPPÉS PAR DES STATISTICIENS TRAVAILLANT À AT&T BELL LABS.

R EST MAINTENANT SOUTENU PAR LA COMMUNAUTÉ ACADÉMIQUE EN SCIENCES INFORMATIQUES ET STATISTIQUES PARTOUT DANS LE MONDE. MINITAB EXPRESS ET R FONCTIONNENT TOUS DEUX SOUS PC OU MAC.



VOS PRÉFÉRENCES INDIVIDUELLES OU VOS GOÛTS VOUS AMÈNERONT À CHOISIR L'UN OU L'AUTRE DE CES LOGICIELS. POUR CERTAINS, RIEN NE PEUT BATTRE LE « GRATUIT ET SOPHISTIQUE » (C'EST R), POUR D'AUTRES, IL FAUT DU « SIMPLE, DU CLAIR ET DU TRANSPARENT » (C'EST MINITAB).

Bien d'autres logiciels de statistiques existent. La plupart ont migré ou vont migrer vers un calcul dans le « CLOUD », de façon à éliminer les problèmes techniques et ceux de compatibilité pour les étudiants comme pour les analystes professionnels de données.

LES GRANDS SYSTÈMES DE BIG DATA DE QUALITÉ INDUSTRIELLE COMME SAS OU JUMP DE SAS INC., SPSS D'IBM, ET STATA SONT LARGEMENT UTILISÉS DANS LES ENTREPRISES ET LA RECHERCHE. ON LES TROUVE TOUS À DES PRIX RÉDUITS POUR LES ÉTUDIANTS. UN ÉTUDIANT AVISÉ DEVRAIT EN PROFITER POUR APPRENDRE L'UN DE CES SYSTÈMES D'ANALYSE DE QUALITÉ INDUSTRIELLE.

POUR CHAQUE PROGRAMME, IL Y A DE NOMBREUX LIVRES ET MANUELS DÉTAILLÉS. UNE BONNE TRENTAINE DE LIVRES OU DE MANUELS SE PRÉTENDENT « UNE INTRODUCTION À R ». NOTRE LIVRE FAVORI DE PROGRAMMATION STATISTIQUE EST *THE LITTLE SAS BOOK* DE L. DELWICHE ET S. SLAUGHTER. ON PEUT Y AJOUTER *THE MINITAB STUDENT HANDBOOK* DE B. RYAN, B. JOINER ET J. CRYER.

LES PACKS STATISTIQUES SONT DIFFÉRENTS EN CE QUI CONCERNE DES DÉTAILS PARFOIS IMPORTANTS. VOUS DEVEZ ÊTRE UN ACHETEUR ÉCLAIRÉ. NOUS VOUS RECOMMANDONS DE CHOISIR UN SYSTÈME QUE VOS COLLÈGUES ONT DÉJÀ TESTÉ. PEU D'ENTRE NOUS SONT TAILLÉS POUR ÊTRE DES PIONNIERS DES LOGICIELS. PENDANT L'APPRENTISSAGE D'UN LOGICIEL, UTILISEZ DES ENSEMBLES PETITS DE DONNÉES FAMILIÈRES. LA PARTIE LA PLUS CÔUTEUSE D'UN LOGICIEL EST LE TEMPS QUE VOUS LUI CONSACREREZ. LA RÈGLE DE BASE POUR APPRENDRE UN LOGICIEL DE STATISTIQUES EST SIMPLE : LA CONNAISSANCE PRÉCÈDE LES RÉSULTATS.

CHERCHER À APPRENDRE SIMULTANÉMENT LA THÉORIE STATISTIQUE ET LA PROGRAMMATION STATISTIQUE EST UN PEU COMME ESSAYER DE MARCHER EN MÂCHANT UN CHEWING-GUM. DIFFÉRENTES APTITUDES ET MÉTHODES DE RÉFLEXION SONT CONCERNÉES EN MÊME TEMPS. PRÉVOYEZ DES TEMPS DIFFÉRENTS POUR APPRENDRE CHAQUE THÈME. ENSUITE COMBINEZ L'ÉTUDE DES DEUX À LA FOIS. DE CETTE FAÇON, VOUS POURREZ DEVENIR UN MÂCHEUR, MARCHEUR, PROGRAMMEUR STATISTICIEN !



Index

A

ABSCISSES, AXE DES, 80
AIRE SOUS LA COURBE, 64-66
AJUSTEMENT EN ANALYSE DE RÉGRESSION, 189-196
ANALYSE DE DONNÉES, 4
ANALYSE DE PUISSANCE, 154-155
ANALYSE DE RÉGRESSION,
 BLUE EN ANALYSE DE RÉGRESSION, 201-202
 CORRÉLATION AU CARRÉ, 196
 DIAGNOSTIC DE RÉGRESSION, 209
 DONNÉES ARRANGÉES, 189, 192, 194-195, 205-207
 ERREUR-TYPE, 203
 ET COEFFICIENTS DE CORRÉLATION, 196
 ET INTERVALLES DE CONFIANCE, 203-206
 ET RÉGRESSION LINÉAIRE, 189-190, 208
 ET VARIABLE ALÉATOIRE DÉPENDANTE, 199-209
 ET VARIABLE DÉPENDANTE, 189
 EXPÉRIENCE DES POIDS D'ÉTUDIANTS, 188-209
 FLUCTUATION D'ERREURS ALÉATOIRES, 199-209
 INFÉRENCE STATISTIQUE, 199-209
 MOYENNE PRÉVUE, 204-206
 MULTIPLE, 208
 PROCESSUS D'AJUSTEMENT, 189-196
 SOMME DES CARRÉS DE LA RÉGRESSION (SCR), 194-196
 SOMME DES CARRÉS DES ERREURS (SCE), 190-195
 TESTS D'HYPOTHÈSES, 207
 VARIABILITÉ DE DONNÉES, 190-195
 VARIABLE DE RÉPONSE, 189
 VARIABLES INDÉPENDANTES, 189, 199-209
ANALYSE DE SÉRIE TEMPORELLE, 214-215
ANALYSE DE VARIANCE, *VOIR* ANOVA
ANALYSE D'IMAGE, 215
ANALYSE DISCRIMINANTE, 213
ANALYSE FACTORIELLE, 213
ANALYSE STATISTIQUE DE DONNÉES MULTIVARIÉES, 212-213
ANALYSE TYPOLOGIQUE, 212
ANOVA (ANALYSE DE VARIANCE), 186, 193-195
 TABLE, 194
APPROXIMATION,
 BINOMIALE, 79-81, 86-88
 CONTINUE, 87-88
 NORMALE, 87-88
ARRONDI, 9
ASPIRINE,
 ESSAIS CLINIQUES, 160-167
 VOIR AUSSI COMPARAISON DE DEUX POPULATIONS
ASTRAGALES, 28
AXE VERTICAL, II

B

BAYES, JOE, 46-50
BAYES, THOMAS, 46-50
BAYES, THÉORÈME DE, 46-50
 VOIR AUSSI FAUX POSITIFS
BAYÉSIEEN, 35
BERNOULLI, JAMES, 79
BERNOULLI, SCHÉMA DE, 74-75, 78
 ET TAILLE D'ÉCHANTILLON, 98-100
BLUE, EN ANALYSE DE RÉGRESSION, 201-202
BIAIS,
 DANS LES ÉCHANTILLONS ALÉATOIRES SIMPLES,
 POUR LES ÉLIMINER, 167
 DANS LES ÉLECTIONS, 126-127
 RÉDUIRE LES BIAIS NATURELS, AVEC COMPARAISONS
 APPARIÉES, 178
BIAIS NATUREL, RÉDUCTION AVEC COMPARAISON APPARIÉE, 178
BINOMIALE, APPROXIMATION, 79-81, 86-88
BLOCS ALÉATOIRES COMPLETS, 184-185
BLOCS EN MÉTHODE EXPÉRIMENTALE, 183-184
BOÎTE À PATTES, 21
BOOTSTRAPPING, 215-216

C

CALCUL DE PROBABILITÉ ET INTERVALLE DE CONFIANCE, 117-119
CALCUL DE Z-SCORE, 84-88, 117-118
CAMÉLÉON AUTOMOBILES,
 COMPARAISON DE MOYENNE DE PETITS ÉCHANTILLONS, 170-171
 INTERVALLE DE CONFIANCE, 134-135
 TEST D'HYPOTHÈSE, 149-150
CARACTÉRISTIQUES
 D'ÉCHANTILLON, 59
 DU MODÈLE, 59
CARRÉ LATIN DANS LES MÉTHODES EXPÉRIMENTALES, 185
CARRÉS DE LA RÉGRESSION, SOMME DES CARRÉS
DE LA RÉGRESSION (SCR), 194-196
CARRÉS DES ÉCARTS, 22, 61-62
CARRÉS DES RÉSIDUS, SOMME DES CARRÉS DES ERREURS
(SCE), 190-195
CHALLENGER (NAVETTE), 3
CHERNOFF, HERMAN, 212
CLAUDE I^{ER}, 28
CLOUS EN LAITON, 98-103
COEFFICIENT,
 BINOMIAL, 76
 DE CORRÉLATION EN ANALYSE DE RÉGRESSION, 196
 DE RÉGRESSION, 191-192
 D'ÉCHANTILLONS, 191-192

- ET LOI MULTIPLICATIVE, 76
- ET TRIANGLE DE PASCAL, 77
- COMMUNICATION, 218
- COMPARAISON,
 - DE MOYENNE DE PETITS ÉCHANTILLONS, 170-171
 - DE SALAIRES MOYENS, 168-169
 - DE TAUX DE SUCCÈS, 160-163
 - DE TAUX D'ÉCHECS, 160-163
 - VOIR AUSSI COMPARAISON DE DEUX POPULATIONS
- COMPARAISON DE DEUX POPULATIONS, 157-179
 - DISTRIBUTION D'ÉCHANTILLONNAGE POUR PROPORTION, 163
 - ET INTERVALLE DE CONFIANCE, 164, 169
 - ET TESTS D'HYPOTHÈSES, 165-167, 169
 - MODÈLE DE, 162
 - MOYENNE DE, 168-169
 - TAUX DE SUCCÈS, 160-163
- COMPARAISONS APPARIÉES, 174-178
 - D'ESSENCES, 174-178
 - ET ÉCARTS-TYPES, 175-176
 - ET TEST t POUR PETITS ÉCHANTILLONS, 176
 - MOYENNE DES, 175-176
- CONFIANCE, INTERVALLE DE, VOIR INTERVALLE DE CONFIANCE
- CONTRÔLE LOCAL DANS LES MÉTHODES EXPÉRIMENTALES, 183
- CONTRÔLE PAR ÉCHANTILLONNAGE, 146-148
- CORRECTION DE CONTINUITÉ, 87-88

D

- DÉS, 28-45
 - PIPÉS, 33
- DEGRÉS DE LIBERTÉ, 131-135
 - ET TESTS D'HYPOTHÈSES, 149-150
 - POUR COMPARAISON DE MOYENNES DE PETITS ÉCHANTILLONS, 171
- DENSITÉ DE PROBABILITÉ, 66
 - DE VARIABLE ALÉATOIRE CONTINUE, 65
- DENSITÉS CONTINUES, PROPRIÉTÉS DES, 66-67
- DIAGRAMME,
 - BRANCHE ET FEUILLE, 12, 18
 - DE POINTS, 9
 - EN BARRE, 11
- DISPERSION,
 - DE PROBABILITÉS, 67
 - DES DONNÉES EN ANALYSE DE RÉGRESSION, 190-192
 - MESURES DE, 19-25
 - VARIANCE DE, 22-23
- DISTRIBUTION BINOMIALE, 77, 81, 83, 86, 88
 - ASYMÉTRIQUE, 82
 - CALCUL POUR VALEUR ÉLEVÉE, 79-80
 - ET FONCTION DE DENSITÉ, 79-80
 - ET LOI NORMALE STANDARD, 82
 - MOYENNE DE, 78
 - VARIANCE DE, 78
- DISTRIBUTION D'ÉCHANTILLONNAGE,
 - DE LA MOYENNE, 104-106
 - DE PROPORTION DE SUCCÈS, 163

- DISTRIBUTION DE PROBABILITÉ,
 - BINOMIALE, 77-78
 - CARACTÉRISTIQUE D'UNE, 59
 - DE VARIABLE ALÉATOIRE, 55-58
 - ET TABLE DE LOIS, 84-85
 - GRAPHIQUE, 56-58
 - MOYENNE DE, 60-61
- DISTRIBUTION NORMALE STANDARD, 79-85
 - ET TABLE DE LOI, 84-85
 - RÈGLE DE CALCUL, 85
- DISTRIBUTION t , 107-109
 - EN COMPARANT DES MOYENNES DE PETITS ÉCHANTILLONS, 171
 - ET INTERVALLE DE CONFIANCE, 131-136
 - ET TESTS D'HYPOTHÈSES, 149-150
 - VALEURS CRITIQUES, 150
- DONNÉES,
 - APPARIÉES OU NON, 177-178
 - ARRANGÉES EN ANALYSE DE RÉGRESSION, 189, 192, 194-195, 205-207
 - CARACTÉRISTIQUES DES, 59
 - ÉCARTS DES, EN ANALYSE DE RÉGRESSION, 190-195
 - MÉDIANE DE, 17
 - MILIEU, 17
 - MOYENNE, 17
 - NOMBRE DE, 11-12, 14-15,
 - TRIÉES, 17
- DONNÉES MULTIVARIÉES,
 - ANALYSE DISCRIMINANTE, 213
 - ANALYSE FACTORIELLE, 213
 - ANALYSE STATISTIQUE, 212-213
 - ANALYSE TYPOLOGIQUE, 212
- DROITE,
 - DE RÉGRESSION AFFINE, 189-192
 - DES MOINDRES CARRÉS, 190

E

- ÉCARTS MOYENS, CARRÉS DES, 22
- ÉCARTS-TYPES,
 - DANS DES COMPARAISONS APPARIÉES, 175-176
 - DANS LES INTERVALLES DE CONFIANCE, 117, 128-130
 - DE POPULATION, 59, 62, 80
 - DE VALEURS MOYENNES, 22, 24-25, 168, 171
 - DÉFINIS PAR RACINE CARRÉE, 23
 - EN COMPARANT DES MOYENNES DE DEUX POPULATIONS, 168
 - EN COMPARANT DES MOYENNES DE PETITS ÉCHANTILLONS, 171
 - ET ÉCHANTILLONNAGE, 101-103, 107
 - ET MESURE DE DISPERSION, 23
 - ET Z-SCORES, 24-25
- ÉCHANTILLON ALÉATOIRE SIMPLE, 92-96
 - VOIR AUSSI ÉCHANTILLONNAGE ALÉATOIRE
- ÉCHANTILLON D'OPPORTUNITÉ, 97
- ÉCHANTILLONNAGE, 89-109
 - D'ACCEPTATION, 150

- EN GRAPPES, 95
- ET ÉCARTS-TYPES, 101-103
- ET EXPÉRIENCE ALÉATOIRE, 98-100, 104-105
- ET INDÉPENDANCE, 93-94, 96
- ET VARIABLES ALÉATOIRES, 98-100, 104-105
- POUR ACCEPTATION, 150
- POUR ÉLIMINER LES BIAIS, 167
- STRATIFIÉ, 95
- SYSTÉMATIQUE, 96-97
- VOIR AUSSI MÉTHODES D'ÉCHANTILLONNAGE
- ÉCHANTILLONNAGE ALÉATOIRE, 95
 - ET ÉLIMINATION DE BIAIS, 167
 - ET INDÉPENDANCE, 92-94, 96
 - UTILISÉ POUR DES INTERVALLES DE CONFIANCE, 114-115, 119
 - SIMPLE, 92-96, 167
- ERREURS,
 - DE TYPE I, 151-154
 - DE TYPE II, 151-154
 - FLUCTUATION D'ERREURS ALÉATOIRES, 199-209
 - HÉTÉROSCÉDASTIQUES, 209
 - MARGE D'ERREUR ET INTERVALLE DE CONFIANCE, 119, 121
 - MESURE D'ERREUR ET DISPOSITIF EXPÉRIMENTAL, 183
 - SOMME DES CARRÉS DES ERREURS (SCE), 190-195
- ERREURS-TYPES,
 - DANS L'ANALYSE DE RÉGRESSION, 203
 - DANS LES INTERVALLES DE CONFIANCE, 118, 128-130
 - EN COMPARANT DES MOYENNES DE DEUX POPULATIONS, 168
 - EN COMPARANT DES MOYENNES DE PETITS ÉCHANTILLONS, 171
 - ET TAILLE D'ÉCHANTILLON, 98-103
- ESPACE ÉCHANTILLON, 30-31, 33, 41
- ESSENCE,
 - COMPARAISONS D', 172-173
 - ET COMPARAISON APPARIÉE, 174-178
 - ET DISPOSITIFS EXPÉRIMENTAUX, 182-186
- ESTIMATEURS, 102-103
 - BLUE EN ANALYSE DE RÉGRESSION, 201-202
 - POUR COMPARER LES MOYENNES DE POPULATIONS, 168-169
- ESTIMATIONS, 102-103, 107
 - D'INTERVALLES DE CONFIANCE, 114-127
- ÉTENDUE INTERQUARTILE, 20-21
- ÉVÉNEMENTS,
 - MUTUELLEMENT EXCLUSIFS, 39, 42, 44
 - PROBABILITÉ D', 35-37
 - RÈGLE D'ADDITION POUR LES, 38-39, 42, 44
 - RÈGLE DE SOUSTRACTION POUR LES, 39, 44
 - RÈGLES DE RÉSULTATS D', 38-39
 - RÉPÉTABLES ET PROBABILITÉS, 35
- EXPÉRIENCE,
 - ALÉATOIRE, 30, 32, 34, 36
 - DES POIDS, 9-12, 16, 18-26
 - ET ÉCHANTILLONS, 98-100, 104-105
 - RÉGRESSION, 188-209
 - DE POIDS, ÉTUDIANTS DE PENN STATE, 9-12, 16, 18-26, 188-209

F

- FAUX POSITIFS, 46-50
- FERMAT, PIERRE DE, 28-45
- FISHER, R. A., 217
- FLUCTUATION D'ERREURS ALÉATOIRES DANS L'ANALYSE DE RÉGRESSION, 199-209
- FONCTION DE DENSITÉ CONTINUE ET DISTRIBUTION BINOMIALE, 79-80
- FRACTILES, 132
- FRÉQUENCES, 10-11, 35, 57-58, 60

G

- GALLUP, 126-127
 - ET TESTS D'HYPOTHÈSES, 143-145
 - NON-PARTICIPATION AU VOTE, 126-127
- GÉNÉRATEUR DE NOMBRES ALÉATOIRES, 65, 94
- GOSSET, WILLIAM, 108-109, 131-132
- GRANDES VALEURS POUR CALCULER DES DISTRIBUTIONS BINOMIALES, 79-80
- GRANDS ÉCHANTILLONS,
 - ET TESTS DE PROPORTIONS, 143-145
 - TESTS D'HYPOTHÈSES DE MOYENNES, 146-148
- GRAPHIQUE,
 - DE DONNÉES, 212
 - DE POINTS VOIR DIAGRAMME
 - DISTRIBUTION DE PROBABILITÉ, 56-58
 - EN BARRE VOIR DIAGRAMME
 - VOIR AUSSI HISTOGRAMMES

H

- HÉTÉROSCÉDASTIQUES, ERREURS, 209
- HISTOGRAMMES,
 - DES FRÉQUENCES, 11, 57-58
 - EFFECTIFS, 57
 - ET ÉCARTS DE MESURE, 19
 - PROBABILITÉ, 56-58
 - SYMÉTRIQUES, 24-25, 77
- HITE, SHERE, 97
- HOLMES, SHERLOCK, 113-130
- HYPOTHÈSE,
 - ALTERNATIVE, 140-141, 152-153, 165-166
 - NULLE (H_0), 140-141, 144-145, 147-150, 152-153, 165-166
 - TEST BILATÉRAL, 144-145
 - TEST UNILATÉRAL DROIT, 144-145
 - TEST UNILATÉRAL GAUCHE, 144-145
 - VOIR AUSSI TESTS D'HYPOTHÈSES
- HYPOTHÈSE ALTERNATIVE (H_A),
 - TEST BILATÉRAL, 144-145
 - TEST UNILATÉRAL DROIT, 144-145
 - TEST UNILATÉRAL GAUCHE, 144-145
- HYPOTHÈSE NULLE (H_0), VOIR HYPOTHÈSE; TESTS D'HYPOTHÈSES

I

IGUANE AUTOS, 170-171
INCERTITUDE, 2
INDÉPENDANCE, 71, 74
 ET ÉCHANTILLON ALÉATOIRE SIMPLE, 92-94, 96
 ET RÈGLE DE MULTIPLICATION SPÉCIALE, 43-44
INFÉRENCE STATISTIQUE, 4
 DANS ANALYSE DE RÉGRESSION, 199-209
INNOVATION, 218
INTÉGRALES, 66-67
INTERVALLES DE CONFIANCE, 112-136
 AUGMENTER LE NIVEAU DES, 121-125
 EN COMPARAISON PAR PAIRE, 176
 ESTIMATION D', 114-127
 ET ANALYSE DE RÉGRESSION, 203-206
 ET CALCUL DE PROBABILITÉ, 117-119
 ET ÉCARTS-TYPES, 117, 128-130
 ET ÉCHANTILLON ALÉATOIRE, 114-115, 119
 ET ERREURS-TYPES, 118, 128-130
 ET MARGE D'ERREUR, 119, 121
 ET MOYENNE DE POPULATION, 128-130, 169
 ET MOYENNES D'ÉCHANTILLON, 130, 171
 ET NIVEAU D'ERREURS, 124-127
 ET PROPORTION DE POPULATION, 128-130
 ET T DE STUDENT, 131-136
 ET TABLE DE FRACTILES, 123
 ET TABLE DE LOI NORMALE, 122
 POUR TAUX DE SUCCÈS, 164
 SIMULATION PAR ORDINATEUR, 120

J-L

JACKKNIFE, 215
JEUX D'ARGENT, 27-45
LANCER DE PIÈCE, 32, 54-55, 58, 60-62, 68-70
LEÇON DE TIR À L'ARC ET INTERVALLE DE CONFIANCE, 116-124

M

MARCHE ALÉATOIRE, 214
MARGE D'ERREUR, VOIR INTERVALLES DE CONFIANCE
MATRICES, 14-15
MÉCANISMES INDÉPENDANTS, 71
MÉDIANE, 17-18, 20-21
MÉRÉ, CHEVALIER DE, 28-29, 75, 78
MESURE,
 DE DISPERSION, 19-25
 D'ERREURS DANS MÉTHODES EXPÉRIMENTALES, 183
 DE VARIABILITÉ, 19-25
MÉTHODE D'ÉCHANTILLONNAGE, 92-94
 ALÉATOIRE, 92-94
 ALÉATOIRE SIMPLE, 92-94
 D'OPPORTUNITÉ, 97
 PAR GRAPPES, 95
 STRATIFIÉ, 95

SYSTÉMATIQUE, 96-97
VOIR AUSSI ÉCHANTILLONNAGE
MÉTHODE EXPÉRIMENTALE,
 BLOCS, 183-184
 CARRÉ LATIN, 185
 CONTRÔLE LOCAL, 183
 D'AJUSTEMENT, 189
 ÉLÉMENTS, 182-183
 ET RANDOMISATION, 183, 185
 ET RÉPLICATION, 183, 185
 ET VARIABILITÉ TOTALE, 186
 MESURE D'ERREUR, 183
 PRINCIPES DE BASES, 183
 TABLEAU 4 X 4, 185
 VARIABILITÉ NATURELLE, 183-185
MODÈLES,
 ALÉATOIRES STOCHASTIQUES, 116-118
 DE RÉGRESSION, 199-202
 POUR DEUX POPULATIONS, 162
MOIVRE, ABRAHAM DE, 79-83, 86-88, 101
MOYENNE, 15-17
 DE DISTRIBUTION BINOMIALE, 78
 DE DISTRIBUTION DE PROBABILITÉ, 60-61
 DE POPULATION, 59, 62
 DE RÉPONSE PRÉVUE DANS ANALYSE DE RÉGRESSION,
 204-206
 DE VARIABLE ALÉATOIRE, 61, 67-69
 D'ÉCHANTILLON ET DISTRIBUTION, 104-106, 171
 D'ÉCHANTILLON ET INTERVALLE DE CONFIANCE, 130, 171
 D'ÉCHANTILLON ET TESTS D'HYPOTHÈSES, 146-148
 DES CARRÉS DES ÉCARTS, 22
 ÉCARTS À LA, 22, 24-25, 168, 171
 ET COMPARAISON APPARIÉE, 175-176
 ET COMPARAISON DE PETITS ÉCHANTILLONS, 170-171
 ET ÉCARTS-TYPES, 22, 24-25, 62, 168, 171
 ET INTERVALLE DE CONFIANCE, 128-130, 169, 171
 ET TESTS D'HYPOTHÈSES, 146-148
 ET TESTS SUR GRANDS ÉCHANTILLONS, 146-148
MOYENNE D'ÉCHANTILLONS,
 COMPARAISON, 170-171
 DISTRIBUTION DE, 104-106
 INTERVALLE DE CONFIANCE, 130, 171
 TESTS D'HYPOTHÈSES, 146-148

N

NIGHTINGALE, FLORENCE, 13
NOMBRES PSEUDO-ALÉATOIRES, 65
NUAGES DE POINTS, 188-189
 ALÉATOIRES, 209

O

OBJECTIVISTE, 35
OPÉRATIONS LOGIQUES, 37
ORDONNÉES, AXE DES, 80

P

PARIS, VOIR JEUX D'ARGENT
PASCAL, BLAISE, 29
PASCAL, TRIANGLE DE, 77
POIDS NUMÉRIQUE, 32
POINTEUR D'UNE ROUE, 63-64
POPULATION,
 CARACTÉRISTIQUES DE LA, 59
 ÉCARTS-TYPES, 59, 62, 80
 ET INTERVALLE DE CONFIANCE, 128-130, 169
 ET TESTS D'HYPOTHÈSES, 146-169
 MOYENNE DE LA, 59, 62
 PROPORTION DE LA, 128-130
PROBABILITÉ CONDITIONNELLE, 40-41
 ET LE PARADOXE DES FAUX POSITIFS, 46-50
 ET LOI DE MULTIPLICATION, 42-44
 D'ÉCHANTILLONNAGE, 100
 D'ERREUR DE TYPE II, 151-155
 NULLE, 63-64
PROBABILITÉS, 4, 27-51
 APPROXIMATION DES, 60
 CARACTÉRISTIQUES DES, 34
 CLASSIQUES, 35
 CONTINUES, 64
 DISCRÈTES, 64, 66
 DISPERSION MOYENNE DES, 67
 ET DISTRIBUTION CUMULATIVE, 84
 ET ÉCHANTILLON, 100
 ET ÉVÉNEMENTS RÉPÉTABLES, 35
 ET FAUX POSITIFS, 46-50
 ET FORMULES, 37-39
 ET RÈGLE DE MULTIPLICATION, 42-44
 NORMALES, 83-85
 PERSONNELLES, 35
 POSITIVES, 34
 TENDANCES CENTRALES DES, 67
PROGRAMME DE CONTRÔLE,
 ANALYSE DE PUISSANCE, 154-155
 PROBABILITÉ D'ERREUR DE TYPE II, 151-155
PROPOSITION CATÉGORIQUE, 2

Q-R

QUALITÉ DE DONNÉES, 217
RACINE CARRÉE, ÉCART-TYPE, 23
RAISONNEMENT DÉDUCTIF, 113
RANDOMISATION, 215-216
 DANS LES MÉTHODES EXPÉRIMENTALES, 183, 185
RÉ-ÉCHANTILLONNAGE, 215-216
RÈGLE D'ADDITION,
 POUR TOUT ÉVÉNEMENT, 38-39, 44
 SPÉCIALE POUR DES ÉVÉNEMENTS MUTUELLEMENT
 EXCLUSIFS, 39, 42, 44
RÈGLE DE MULTIPLICATION,
 ET COEFFICIENT BINOMIAUX, 76

 ET PROBABILITÉ CONDITIONNELLE, 42-44
 SPÉCIALE ET INDÉPENDANCE, 43-44
 SPÉCIALE ET PROBABILITÉ CONDITIONNELLE, 42-44
RÈGLE DE SOUSTRACTION POUR DES ÉVÉNEMENTS, 39, 44
RÉGRESSION, 187-209
 LINÉAIRE, 189-190, 208
 NON LINÉAIRE DANS L'ANALYSE DE RÉGRESSION, 208
RÉPLICATION DANS LES MÉTHODES EXPÉRIMENTALES, 35
REPRÉSENTATION GRAPHIQUE, 13
RÉSULTATS,
 D'ÉVÉNEMENTS, 38-39
 ÉLÉMENTAIRES, 30, 32-38, 41
 NUMÉRIQUES ET ÉCHANTILLONS, 98-100, 104-105
RÉSUMÉ DE DONNÉES, 12
RÉSUMÉ STATISTIQUE, 14-26
 DANS LES TESTS D'HYPOTHÈSES, 148

S

SALK, VACCIN CONTRE LA POLIO, 3
SCE,
 PAR RAPPORT AUX VARIATIONS DE DONNÉES, 193-195
 SOMME DES CARRÉS DES ERREURS, 190-195
SCR, SOMME DES CARRÉS DE LA RÉGRESSION, 194-196
SÉLECTION ALÉATOIRE DE JURY, BIAIS RACIAUX, 138-141
SEUIL DE CONFIANCE, 122
SEUIL DE SIGNIFICATION,
 ET TESTS D'HYPOTHÈSES, 141-142, 145, 147-148
 ET TRAVAIL SCIENTIFIQUE, 142
 FIXE, 142
SIGMA, 16
SOMMATION, 16
SOMME DES CARRÉS,
 DE LA RÉGRESSION, DANS L'ANALYSE DE RÉGRESSION,
 194-196
 DES ERREURS, DANS L'ANALYSE DE RÉGRESSION, 190-195
 DES ERREURS, RELATIVE AUX ÉCARTS DE DONNÉES, 193-195
SONDAGE,
 ÉLECTION, 114-127
 ET BIAIS, 126-127
 ET NIVEAU D'ERREURS, 124-127
SONDAGE ÉLECTORAL, 114-127
 ET TESTS D'HYPOTHÈSES, 143-145
SONDAGES GALLUP, VOIR GALLUP
STATISTIQUE DE TEST,
 DANS LES TESTS D'HYPOTHÈSES, 140-141, 144-145, 147-148,
 165-166, 169
 ET PETITS ÉCHANTILLONS POUR COMPARAISON
 APPARIÉE, 176
STATISTIQUES,
 DE MORTALITÉ, 13
 DESCRIPTIVES (RÉSUMÉS), 14-26, 148
STUDENT, VOIR *t* DE STUDENT
SUBJECTIVISTE, 35
SUCCÈS, NOMBRE DE, 75

T

TABLE,
 DE DÉCISION, 152
 DE DISTRIBUTION NORMALE, 78
 DE FRACTILES, 132
 FACTORIELLE 2 X 2, 184-185
 TABLES DE FRÉQUENCES, *VOIR* HISTOGRAMMES
 TAILLE D'ÉCHANTILLON, 91
 AUGMENTER LA TAILLE, 124-125
 COMPARAISON DE PETITS ÉCHANTILLONS, 170-171
 ET ERREURS-TYPES, 98-103
 ET GRANDS TAILLES, 143-148
 ET NIVEAU DE CONFIANCE, 124-125
 TAUX DE DÉCÈS, *VOIR* STATISTIQUES
 TAUX DE SUCCÈS, 99
 COMPARAISON DE DEUX POPULATIONS, 160-163
 DANS LES TESTS D'HYPOTHÈSES, 143-145
 ET DISTRIBUTION D'ÉCHANTILLONNAGE, 163
 ET INTERVALLES DE CONFIANCE, 164
 TAUX D'ÉCHECS, COMPARAISON POUR DEUX POPULATIONS, 160-163
 TENDANCE CENTRALE, 14
 MÉDIANE, 17-18
 MOYENNE, 15-16
VOIR AUSSI VARIABILITÉ
 TEST DE SIGNIFICATION,
 POUR CONTRÔLE PAR ÉCHANTILLONNAGE, 146-148
 POUR PROPORTIONS, 143-145
 TESTS D'HYPOTHÈSES, 138-139
 EN ANALYSE DE RÉGRESSION, 207
 EN COMPARAISON PAR PAIRE, 176
 ET DEGRÉS DE LIBERTÉS, 144-150
 ET MOYENNE DE POPULATION, 147-148, 169
 ET SEUIL DE SIGNIFICATION, 141-142, 145
 ET THÉORIE DE LA DÉCISION, 151-155
 GRANDS ÉCHANTILLONS ET MOYENNE DE POPULATION, 146-148
 GRANDS ÉCHANTILLONS ET TEST DE PROPORTIONS, 143-145
 STATISTIQUES, 140-142, 144-145, 147-148, 165-166, 169
 THÉORÈME CENTRAL LIMITE, 106, 169, 83-88
 DÉFAUTS, 107
 THÉORIE DE LA DÉCISION, TEST D'HYPOTHÈSE, 151-155
 TRAITEMENTS EXPÉRIMENTAUX, 182-183
 TRAVAIL D'ÉQUIPE, 218
 TUKEY, JOHN, 12, 21

U

UNILATÉRAL GAUCHE, HYPOTHÈSE ALTERNATIVE, 144-145
 UNITÉS EXPÉRIMENTALES, 182-183

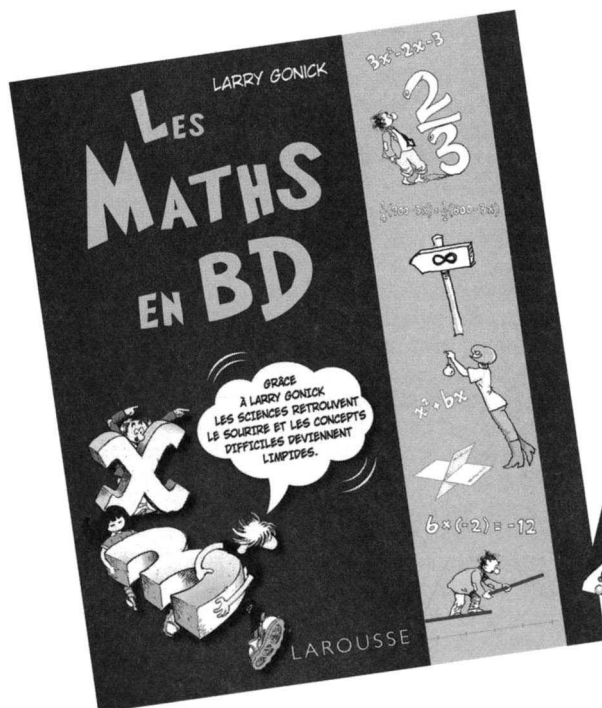
V

VALEUR,
 ESPÉRÉE, 61
 OBSERVÉE DE T , 149-150
 OBSERVÉE DE Z ET TESTS D'HYPOTHÈSES, 144-145, 165-166, 169
 P DANS TESTS D'HYPOTHÈSES, 141-142, 148
 TYPIQUE, 14-18
 VALEURS,
 CRITIQUES, 132-136, 150
 EXTRÊMES, 18, 21-23
 VARIABILITÉ,
 NATURELLE ET MÉTHODE EXPÉRIMENTALE, 183-185
 RÉDUIRE AVEC COMPARAISONS APPARIÉES, 178
 TOTALE DUE À LA RÉGRESSION, 194-195
 TOTALE ET MÉTHODE EXPÉRIMENTALE, 186
 VARIABILITÉ NATURELLE,
 DANS LES MÉTHODES EXPÉRIMENTALES, 183-185
 RÉDUIRE AVEC COMPARAISONS APPARIÉES, 178
 VARIABLE,
 ALÉATOIRE BINOMIALE, 74-76, 140
 ALÉATOIRE CONTINUE, 63, 65, 67
 ALÉATOIRE DÉPENDANTE DANS L'ANALYSE DE RÉGRESSION, 199-209
 ALÉATOIRE DISCRÈTE, 63
 DÉPENDANTE DANS L'ANALYSE DE RÉGRESSION, 189
 INDÉPENDANTE DANS L'ANALYSE DE RÉGRESSION, 189, 199-209
 VARIABLES ALÉATOIRES, 53-72
 DISTRIBUTION DE PROBABILITÉ, 55-58
 ET ÉCHANTILLONNAGE, 98-100, 104-105
 MOYENNE DE, 61, 67-69
 SOMMATION DE, 68-71
 T DE STUDENT, 107-109
 VARIANCE DE, 62, 67-71
 VARIABLES ALÉATOIRES CONTINUES, 63
 ET DENSITÉ DE PROBABILITÉ, 65
 ET MOYENNE, 67
 ET VARIANCE, 67
 VARIABLES DE RÉPONSE DANS ANALYSE DE RÉGRESSION, 189
 VARIANCE,
 ANALYSE DE (ANOVA), *VOIR* ANOVA
 DE DISTRIBUTION BINOMIALE, 78
 DE VARIABLES ALÉATOIRES, 62, 67-71
 DE VARIABLES ALÉATOIRES CONTINUES, 67
 D'ÉCHANTILLON, 22
 DES ÉCARTS, 22-23

Z

Z , VALEUR OBSERVÉE, TESTS D'HYPOTHÈSES, 144-145, 165-166, 169
 Z -SCORE ET ÉCARTS-TYPES, 24-25

D'AUTRES « SCIENCES EN BD »
AVEC LARRY GONICK
CHEZ LAROUSSE...



LES MATHS EN BD

LARRY GONICK

AVEC CE VOLUME SPÉCIALEMENT CONSACRÉ À L'ALGÈBRE, RÉVISEZ TOUS LES FONDAMENTAUX : RAPPELS SUR LES NOMBRES ET LES OPÉRATIONS, MISE EN ÉQUATION DES VARIABLES ET SIMPLIFICATIONS, ÉQUATIONS DU PREMIER ET DU SECOND DEGRÉS, GRAPHIQUES, RACINES CARRÉES... AVEC DES CAS CONCRETS ET DES PROBLÈMES CORRIGÉS POUR S'ENTRAÎNER !

240 PAGES

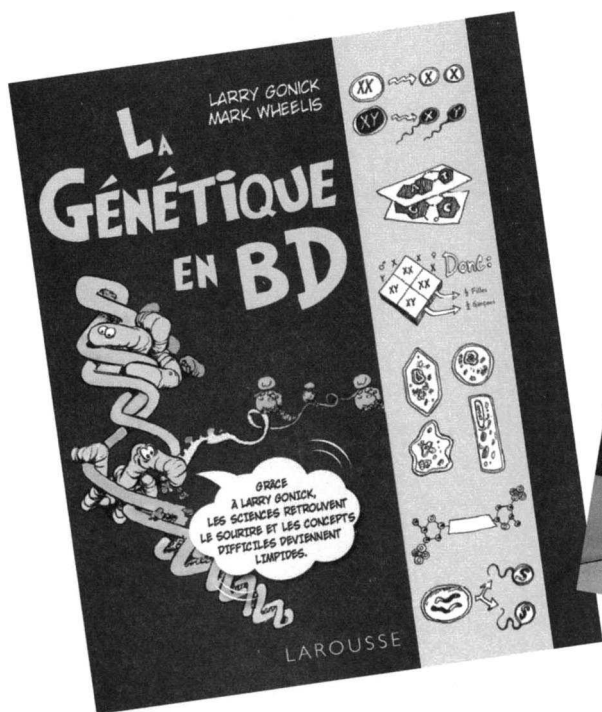


LA CHIMIE EN BD

LARRY GONICK ET CRAIG CRIDDLE

PASSEZ EN REVUE TOUS LES DOMAINES DE LA CHIMIE : LA MATIÈRE ET LES ÉLÉMENTS DU TABLEAU PÉRIODIQUE, LE PROTON, LE NEUTRON ET L'ÉLECTRON (ET SON RÔLE DANS LA LIAISON CHIMIQUE), LES IONS, LES MOLÉCULES, LA MASSE MOLAIRE, LES GAZ PARFAITS, LA THERMODYNAMIQUE... AVEC DES CAS CONCRETS ET DES PROBLÈMES CORRIGÉS POUR S'ENTRAÎNER !

256 PAGES

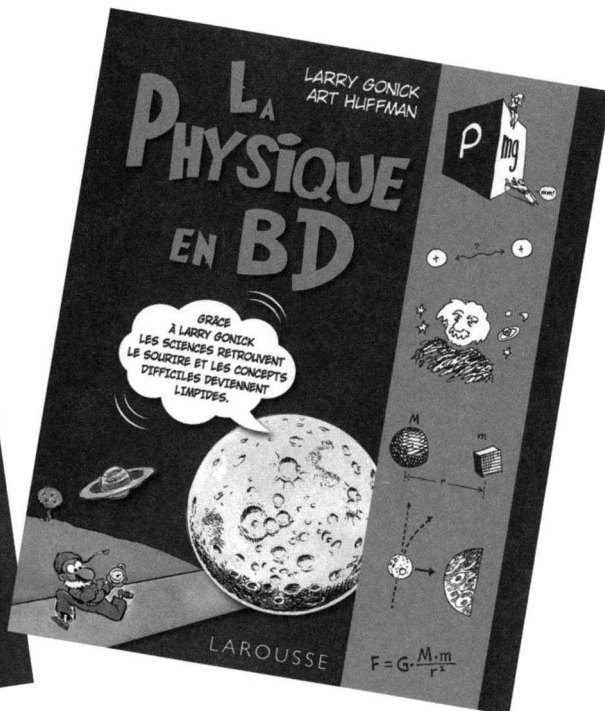


LA GÉNÉTIQUE EN BD

LARRY GONICK ET MARK WHEELIS

DÉCOUVREZ OU RÉVISEZ LES MÉCANISMES INCROYABLES DE LA GÉNÉTIQUE : LA REPRODUCTION ET L'HÉRÉDITÉ, LES LOIS DE MENDEL, DOMINANCE ET RÉCESSIVITÉ, LES CHROMOSOMES, LA DÉCOUVERTE DE L'ADN, LES MUTATIONS, LA SYNTHÈSE DES PROTÉINES, LA RÉGULATION DES GÈNES...

224 PAGES



LA PHYSIQUE EN BD

LARRY GONICK ET ART HUFFMAN

POUR RÉVISER LES FONDAMENTAUX DE LA PHYSIQUE EN S'AMUSANT, UN OUVRAGE EN DEUX PARTIES : MÉCANIQUE (L'ATTRACTION UNIVERSELLE, LES LOIS DE NEWTON, LES FORCES, L'ÉNERGIE CINÉTIQUE, L'ÉNERGIE POTENTIELLE, ETC.) ET ÉLECTRICITÉ (LA CHARGE, LES CHAMPS ÉLECTRIQUES ET MAGNÉTIQUES, LES CAPACITÉS, LE COURANT ÉLECTRIQUE, MONTAGES EN SÉRIE ET EN PARALLÈLE, ETC.).

224 PAGES



LAROUSSE s'engage pour
l'environnement en réduisant
l'empreinte carbone de ses livres.
Celle de cet exemplaire est de :
1,5 kg éq. CO₂
Rendez-vous sur
www.larousse-durable.fr

IMPRIMÉ EN ESPAGNE CHEZ UNIGRAF
DÉPÔT LÉGAL : SEPTEMBRE 2016
318486-01/11033648 – SEPTEMBRE 2016

LES STATISTIQUES EN BD

AUJOURD'HUI, PRESQUE TOUS LES DOMAINES D'ACTIVITÉ UTILISENT LES STATISTIQUES : ÉCONOMIE, FINANCE, TECHNOLOGIE, INFORMATIQUE, BIOLOGIE, COMMUNICATION... AVEC CE VOLUME QUI PASSE EN REVUE LES GRANDS PRINCIPES DE CETTE DISCIPLINE, NOTAMMENT LES PLUS POINTUS, RÉVISEZ VOS FONDAMENTAUX EN VOUS AMUSANT ET RÉVÉLEZ LE PETIT GÉNIE QUI SOMMEILLE EN VOUS.

- Au programme : données, descriptions, probabilités, écarts-types, variables aléatoires, distributions, lois normale et binomiale, échantillonnages, intervalles de confiance, tests d'hypothèse, méthodes expérimentales, régression linéaire...
- Des illustrations qui éclairent les énoncés scientifiques tout en y injectant une bonne dose d'humour. Avec ses mises en situation désopilantes et ses cas pratiques, Larry Gonick parvient à rendre les statistiques modernes, utiles et passionnantes.

LARRY GONICK est un auteur de bandes dessinées et un scientifique américain. Il a enseigné les mathématiques à Harvard et participe au développement du journalisme scientifique au MIT (Massachusetts Institute of Technology). Ses nombreux « comics » de vulgarisation scientifique (chimie, histoire de l'Univers, informatique, physique, génétique...) sont des best-sellers outre-Atlantique, où ils sont devenus des classiques.

WOOLLCOTT SMITH est professeur émérite de statistiques à l'université de Temple à Philadelphie. Il est titulaire d'une licence et d'un master de l'université du Michigan et d'une thèse de l'université Johns-Hopkins.

17,99 € PRIX TTC FRANCE

ISBN 978-2-03-593291-4

3519705



9 782035 932914

Illustrations de couverture :
Larry Gonick

LAROUSSE
DICTIONNAIRES

www.editions-larousse.fr